



Embracing heterogeneity: Coalescing the tree of life and the future of phylogenomics

Downloaded from: <https://research.chalmers.se>, 2023-05-06 02:33 UTC

Citation for the original published paper (version of record):

Bravo, G., Antonelli, A., Bacon, C. et al (2019). Embracing heterogeneity: Coalescing the tree of life and the future of phylogenomics. PeerJ, 2019(2). <http://dx.doi.org/10.7717/peerj.6399>

N.B. When citing this work, cite the original published paper.

Embracing heterogeneity: coalescing the Tree of Life and the future of phylogenomics

Gustavo A. Bravo¹, Alexandre Antonelli^{1,2,3,4}, Christine D. Bacon^{2,3}, Krzysztof Bartoszek⁵, Mozes P. K. Blom⁶, Stella Huynh⁷, Graham Jones³, L. Lacey Knowles⁸, Sangeet Lamichhaney¹, Thomas Marcussen⁹, Hélène Morlon¹⁰, Luay K. Nakhleh¹¹, Bengt Oxelman^{2,3}, Bernard Pfeil³, Alexander Schliep¹², Niklas Wahlberg¹³, Fernanda P. Werneck¹⁴, John Wiedenhoeft^{12,15}, Sandi Willows-Munro¹⁶ and Scott V. Edwards^{1,17}

¹ Department of Organismic and Evolutionary Biology, Museum of Comparative Zoology, Harvard University, Cambridge, MA, USA

² Gothenburg Global Biodiversity Centre, Göteborg, Sweden

³ Department of Biological and Environmental Sciences, University of Gothenburg, Göteborg, Sweden

⁴ Gothenburg Botanical Garden, Göteborg, Sweden

⁵ Department of Computer and Information Science, Linköping University, Linköping, Sweden

⁶ Department of Bioinformatics and Genetics, Swedish Museum of Natural History, Stockholm, Sweden

⁷ Institut de Biologie, Université de Neuchâtel, Neuchâtel, Switzerland

⁸ Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI, USA

⁹ Centre for Ecological and Evolutionary Synthesis, University of Oslo, Oslo, Norway

¹⁰ Institut de Biologie, Ecole Normale Supérieure de Paris, Paris, France

¹¹ Department of Computer Science, Rice University, Houston, TX, USA

¹² Department of Computer Science and Engineering, Chalmers University of Technology and University of Gothenburg, Göteborg, Sweden

¹³ Department of Biology, Lund University, Lund, Sweden

¹⁴ Coordenação de Biodiversidade, Programa de Coleções Científicas Biológicas, Instituto Nacional de Pesquisa da Amazônia, Manaus, AM, Brazil

¹⁵ Department of Computer Science, Rutgers University, Piscataway, NJ, USA

¹⁶ School of Life Sciences, University of Kwazulu-Natal, Pietermaritzburg, South Africa

¹⁷ Gothenburg Centre for Advanced Studies in Science and Technology, Chalmers University of Technology and University of Gothenburg, Göteborg, Sweden

Submitted 5 February 2018

Accepted 7 January 2019

Published 14 February 2019

Corresponding author

Gustavo A. Bravo,
gustavo_bravo@fas.harvard.edu

Academic editor

Chris Creevey

Additional Information and
Declarations can be found on
page 36

DOI 10.7717/peerj.6399

© Copyright
2019 Bravo et al.

Distributed under
Creative Commons CC-BY 4.0

OPEN ACCESS

ABSTRACT

Building the Tree of Life (ToL) is a major challenge of modern biology, requiring advances in cyberinfrastructure, data collection, theory, and more. Here, we argue that phylogenomics stands to benefit by embracing the many heterogeneous genomic signals emerging from the first decade of large-scale phylogenetic analysis spawned by high-throughput sequencing (HTS). Such signals include those most commonly encountered in phylogenomic datasets, such as incomplete lineage sorting, but also those reticulate processes emerging with greater frequency, such as recombination and introgression. Here we focus specifically on how phylogenetic methods can accommodate the heterogeneity incurred by such population genetic processes; we do not discuss phylogenetic methods that ignore such processes, such as concatenation or supermatrix approaches or supertrees. We suggest that methods of data acquisition and the types of markers used in phylogenomics will remain restricted until a posteriori methods of marker choice are

made possible with routine whole-genome sequencing of taxa of interest. We discuss limitations and potential extensions of a model supporting innovation in phylogenomics today, the multispecies coalescent model (MSC). Macroevolutionary models that use phylogenies, such as character mapping, often ignore the heterogeneity on which building phylogenies increasingly rely and suggest that assimilating such heterogeneity is an important goal moving forward. Finally, we argue that an integrative cyberinfrastructure linking all steps of the process of building the ToL, from specimen acquisition in the field to publication and tracking of phylogenomic data, as well as a culture that values contributors at each step, are essential for progress.

Subjects Biodiversity, Computational Biology, Evolutionary Studies, Genomics

Keywords Gene flow, Genome, Multispecies coalescent model, Retroelement, Speciation, Transcriptome

INTRODUCTION

Phylogenomics has been greatly enriched with the introduction of high-throughput sequencing (HTS) and increased breadth of phylogenomic sampling, which have allowed researchers interested in the Tree of Life (ToL) to scale up in several dimensions and placing both fields squarely in the era of ‘big data’. Additionally, conceptual advances and improvements of statistical models used to analyze these data are helping bridge what some have perceived as a gap between phylogenetics and phylogeography (e.g., *Felsenstein, 1988; Huson & Bryant, 2006; Edwards et al., 2016a*). However, as datasets become larger, researchers are inevitably faced with a plethora of heterogeneous signals across different genomic regions that often depart from a dichotomously-branching phylogeny (*Kunin, Goldovsky & Darzentas, 2005; Jeffroy et al., 2006; Mallet, Besansky & Hahn, 2015*). These signals cover an increasingly large array of biological processes at the level of genes and genomes, as well as individual organisms and populations, including processes such as recombination, hybridization, gene flow, and polyploidization; they can be thought of as conflicting, but in truth, they are simply a record of the singular history that we commonly refer to as the ToL (*King & Rokas, 2017*). One of the grand challenges of evolutionary biology is deciphering this history, whether at the level of genes, populations, species, or genomes, however, currently phylogenomicists have not yet determined how to fully exploit the diverse signals in molecular data.

In this perspective piece, our goal is to highlight opportunities for the field of phylogenomics to overcome conceptual and practical challenges as we navigate our way through the era of big data. We argue that conceptually and analytically embracing heterogeneity generated by population-level processes and associated with different histories across the genome will lead to increased insight into the ToL and its underlying processes. Because major theoretical advances in phylogenomics and population genetics resulting from the advent of genome-scale data have been reviewed elsewhere (e.g., *Delsuc, Brinkmann & Philippe, 2005; Degnan & Rosenberg, 2009; Ellegren & Galtier, 2016; Payseur & Rieseberg, 2016; Hahn, 2018*), here we focus on advances and

challenges to phylogenomics specifically brought about by population genetic processes, which inevitably leads us to focus on the major conceptual framework dealing with such processes, the multispecies coalescent model (MSC). Likewise, advances in other methods such as supertrees (reviewed by [Warnow, 2018](#)), supermatrices ([de Queiroz & Gatesy, 2007](#); [Philippe et al., 2017](#)), and partitioned analyses (e.g., [Lanfear et al., 2014](#); [Kainer & Lanfear, 2015](#)) are not central to the objectives of this perspective piece and further details on those topics can be found elsewhere.

A key concept introduced by the scaling up from phylogeography to phylogenomics is the continuum of processes and analytical methods—the so-called phylogeography-phylogenetics continuum ([Edwards et al., 2016a](#)). We argue here that bridging this continuum is critical for advancing phylogenetics. This bridging can be done by either developing phylogenomic approaches that acknowledge and explicitly account for phylogeographic processes, or by determining the regions of parameter space (e.g., branch lengths in tree, level of gene flow) if any, where such within-species processes are not relevant. For example, the choice of markers in a given phylogenomics project is currently guided more by convenience and cost than by evaluating the biological properties and phylogenomic signals in those data; but comparisons of signals across various types of markers (e.g., transcriptomes, noncoding regions) reveal that marker choice is a critical step toward shedding light on the history of populations and unraveling potential processes underlying such history ([Rokas et al., 2003](#); [Cutter, 2013](#); [Jarvis et al., 2014](#); [Reddy et al., 2017](#)). On the analysis side, we are in desperate need of methods that can handle the increasingly large data sets being produced by empiricists, but at the same time there is a desire to include increasingly diverse sources of signal in estimates of divergence times, biogeographic history, and models of diversification ([Delsuc, Brinkmann & Philippe, 2005](#); [Jeffroy et al., 2006](#); [Kumar et al., 2012](#)). Finding the balance between breadth, depth, and computational feasibility in project design and statistical analysis is crucial for the field today.

Additionally, we discuss ways in which data archiving and integration can benefit access to phylogenomic data and the contributions of phylogenomics to society. Are the priorities that society places on the many disciplines feeding into scientific efforts toward the ToL—fieldwork, museum collections, databases—appropriate for this grand mission? Although we cannot possibly answer all of these questions within the scope of this perspective, we hope to at least spur discussion on the wide range of field, laboratory, conceptual, and societal issues that allow phylogenomics to move forward.

We first describe the types of genomic data that researchers can generate to perform phylogenomic analyses and how those are more or less suitable for phylogenomic and phylogeographic analyses. We then discuss key concepts around the MSC and highlight the need to expand this model beyond its current limitations. We then discuss how the interplay between phylogenomics and macroevolution might strengthen our understanding of diversity patterns and offer suggestions as to how the community can overcome limitations posed by current methods and models in both fields. Finally, we discuss desired practices that, as a community, phylogenomicists, museum scientists,

field biologists, bioinformaticians, and other scientists can embrace toward the goal of assembling the ToL.

SURVEY METHODOLOGY

During the ‘Origins of Biodiversity Workshop’ organized during May 15–19, 2017 by Chalmers University of Technology and the University of Gothenburg, Sweden, under the auspices of the Gothenburg Centre for Advanced Studies (GoCAS), we gathered scholars and students from several countries and scientific backgrounds to discuss future perspectives in the fields of phylogenomics and phylogeography. We spent one week sharing our recent experiences in these fields and outlining the topics presented here and continued remotely to complete this review. Our goal is not to provide a complete overview of phylogenomic and phylogeographic research, but rather present a number of conceptual and practical aspects that we feel are essential to keep the momentum that these fields have gained during recent times.

DATA GENERATION AND DATA TYPES IN PHYLOGENOMICS

The need for large-scale phylogenomic data

One of the fundamental challenges in evolutionary biology is to estimate a ToL for all species. The potential impact of such large phylogenies is reflected in their publication in the highest impact journals, but also in their broad contribution, which extends beyond big data, to methodological innovations, and downstream understanding of macroevolutionary processes (e.g., coalescent methods of species tree inference; accounting for hybridization and unsampled species or localities in datasets; understanding community or genome evolution through large-scale phylogenetics). Hence, the phylogenomics community is now placing a high priority on very large-scale trees, whether in terms of number of taxa, number of genes, or both. The current need for large phylogenies and the high priority placed on them by high-impact journals can also result in short-cuts, wherein extant species lacking any molecular data are often placed in trees based on current taxonomy (Jetz *et al.*, 2012; Zanne *et al.*, 2014; Faurby & Svenning, 2015), which can result in conflicts with more robust estimates based on actual data (Brown, Wang & Smith, 2017). At the same time, however, hypothesis-testing in areas such as macroevolution, macroecology, biodiversity, and systematics require these large-scale trees, even as they present challenges being built on high quality data. The phylogenetic knowledge on which we lay a foundation for downstream analyses must be robust, and therefore it is essential that the input phylogenetic hypotheses themselves are robust (Pyron, 2015). Indeed, the current bottlenecks in large-scale phylogenomic data do not appear to be the sequencing, but rather the compilation of high quality, well-curated genomic resources and the availability of adequate software and methods that can fuel phylogenomics for the next century (e.g., Global Genome Initiative, www.mnh.si.edu/ggi/).

Data quality

Genome-scale data in the form of multiple alignments and other homology statements are the foundation of phylogenomics. A major challenge is the difficulty of comprehensive

quality checks of data, given that HTS datasets are so large. As researchers collect datasets consisting of thousands of alignments across scores of species, data quality is a serious concern that is left for detection and handling primarily by computer algorithms. In addition to inherent systematic errors in the data ([Kocot et al., 2017](#)), several examples of errors in phylogenomic data sets have been reported in the literature, including the use of unintended paralogous sequences in alignments (e.g., [Struck, 2013](#)); mistaking the genome sequence of one species for another ([Philippe et al., 2011](#)); and inclusion of genome sequence from parasites into the genome of the host ([Kumar et al., 2013](#)). However, the incidence of smaller errors in alignments that are not easily discerned from natural allelic variation, such as base miscalls or misplaced indels, are probably much more widespread than has been reported in the literature. Combined with the sensitivity of some phylogenomic datasets to individual loci or single nucleotide polymorphisms (SNPs) within loci ([Shen, Hittinger & Rokas, 2017](#)), such errors could have damaging consequences for phylogenomic studies, for the inference of topology and branch lengths of both gene trees and species phylogenies ([Marcussen et al., 2014](#); [Bleidorn, 2017](#)). Furthermore, as phylogenomic datasets increase in size, it is likely that the accumulation of errors due to the use of different sequencing chemistries and sequencing depths ([Quail et al., 2012](#); [Goodwin, McPherson & McCombie, 2016](#)) will ultimately influence phylogenetic inference. We predict that the impact of these errors will largely depend on the sampling breadth and taxonomic scale of each study, and whether the species phylogeny is a tree or a network.

Sequencing high-quality samples from well-archived voucher specimens is a good first step to increase reproducibility and alleviate issues related to sample identity ([Peterson et al., 2007](#); [Pleijel et al., 2008](#); [Chakrabarty, 2010](#); [Turney et al., 2015](#); [Troudet et al., 2018](#)). For individual phylogenomic studies, wholesale manual inspection of every locus is unsustainable ([Irisarri et al., 2017](#)), but spot checks of a subset of the data (e.g., 5–10% of the alignments) is a recommended best practice ([Philippe et al., 2011](#)) that is beginning to be encouraged in peer review and in published papers ([Montague et al., 2014](#); [Liu et al., 2017](#)). Such checking is important not only for new data generated by a given study, but also for data downloaded from public repositories such as NCBI and OrthoMaM ([Ranwez et al., 2007](#); [Douzery et al., 2014](#)), which are well known to contain errors ([Wesche, Gaffney & Keightley, 2009](#)). Because several databases do not include the raw sequence data it is often impossible to evaluate whether oddities may derive from poor sequencing. Robust pipelines for flagging poorly aligned sites or non-homologous sequences, based on existing tools or novel scripts such as Gblocks ([Castresana, 2000](#); [Talavera & Castresana, 2007](#)) or TrimAl ([Capella-Gutierrez, Silla-Martinez & Gabaldon, 2009](#)) are gradually being put into practice ([Marcussen et al., 2014](#); [He et al., 2016](#); [Irisarri et al., 2017](#)).

Coding regions, whether derived from transcriptomes or whole-genome data, are particularly amenable to spot checking of alignments and to filtering out of low-quality data with bioinformatic pipelines (e.g., [Dunn, Howinson & Zapata, 2013](#); [Blom, 2015](#)). Coding regions have the advantage of allowing amino acids to guide alignments, which is particularly useful for highly divergent sequences. Stop codons can help flag errors

or genuine pseudogenes. Examining gene tree topologies is also widely used to detect paralogs in phylogenomic data (e.g., [Betancur, Naylor & Ortí, 2014](#)). Examining gene trees for aberrantly long branch lengths can also reveal misalignments (e.g., [He et al., 2016](#)); sensitivity analyses of methods for indirectly detecting errors in alignments are sorely needed.

Data generation and marker development

Genome reduction methods

A growing number of genome reduction methods are now providing empiricists with the means to generate genomic subsets suitable for phylogenetic and phylogeographic inference (reviewed by [McCormack et al., 2013](#); [Leaché & Oaks, 2017](#); [Lemmon & Lemmon, 2013](#)). For phylogenomics, most prominently featured are sequence-capture, focusing on highly conserved regions (e.g., [Faircloth et al., 2012](#); [Lemmon, Emme & Lemmon, 2012](#); reviewed by [Jones & Good, 2016](#)) and transcriptomes (e.g., [Misof et al., 2014](#); [Cohen et al., 2016](#); [Fernández et al., 2014](#); [Park et al., 2015](#); [Simion et al., 2017](#); [Irisarri et al., 2017](#)), but phylogenomic trees have also been constructed based on restriction-digest methods that primarily focus on SNPs ([Leaché, Chavez & Jones, 2015](#); [Harvey et al., 2016](#)) and analysis of transposable elements (e.g., [Suh, Smeds & Ellegren, 2015](#)). This diversity of marker types for phylogenetics should be celebrated, but each marker type brings with it a list of pros and cons. For example, many questions in the higher level phylogenetics of animals and plants have so far relied almost exclusively on transcriptome data. However, the uncritical use of transcriptomes in phylogenetics is not without caveats. At high taxonomic levels, coding regions can exhibit extreme levels of among-taxon base composition, sometimes resulting in strong violations of phylogenetic models ([Romiguier et al., 2016](#); [Romiguier & Roux, 2017](#)). Coding regions can exhibit reduced levels of incomplete lineage sorting (ILS) compared to noncoding regions ([Scally et al., 2012](#)). Such reduced ILS could in fact be helpful in building complex phylogenies with rapid radiations ([Edwards, 2009a](#)), but it will certainly distort estimated branch lengths when coalescent methods, which assume neutrality, are used. In addition to yielding abundant sequence variations and SNPs, transcriptome data also yields information on the levels of expression of various genes, and in which tissue-specific genes are expressed. However, using these aspects of transcriptome data are less likely to be of use to phylogeneticists, precisely because specific genes are often tissue-specific and because expression data can exhibit high levels of variation among individuals, populations and species in space and time ([Todd, Black & Gemmell, 2016](#)). Such variation will certainly pose limitations for long term phylogenomic endeavors that will likely combine data collected originally for different purposes.

Although non-coding portions of the genome have been largely neglected in phylogenomics because they are difficult to align and analyze, we are now making progress in understanding their sequence evolution and how it might be informative for comparative purposes ([Ulitsky, 2016](#); [Edwards, Cloutier & Baker, 2017](#)). For instance, studies on enhancer and promoter evolution in mammals have shown that despite low levels of sequence conservation, there is conservation of regulatory function

and 3D structure across species that carries information for comparative purposes (Villar *et al.*, 2015; Berthelot *et al.*, 2018). The development of methods to infer and interpret the evolutionary history and phylogenetic signal of non-coding elements and 3D genome structure is a critical priority.

Single nucleotide polymorphisms have been advocated by some authors for higher level phylogenetics (Leaché & Oaks, 2017), but the available methods for analyzing such data are still extremely limited. For example, concatenation and two coalescent methods (SNAPP and SVD quartets: Bryant *et al.*, 2012; Chifman & Kubatko, 2014) have recently been highlighted as the main methods available for phylogenomic analysis of SNPs (Leaché & Oaks, 2017). But each of these methods has its shortcomings. It is likely that concatenation of SNPs will be misleading for many of the same reasons that concatenating genes can be misleading, due to different gene histories (see section ‘Concepts and models in phylogenomics’ for further details; Kubatko & Degnan, 2007). SNAPP, a coalescent method suitable for analysis of SNPs (Bryant *et al.*, 2012), works well only on relatively small data sets, and it is unclear how well SVD quartets performs on some data sets (Shi & Yang, 2018). Although SNPs do provide a helpful route around the often-violated assumption in coalescent models of no recombination within loci (Bryant *et al.*, 2012), and are informative markers for phylogeography and population genetics (Brumfield *et al.*, 2003), it remains to be seen how powerful they are at higher phylogenetic levels.

Despite the diversity of marker types for phylogenomics, it remains unclear whether features specific to each marker type can ultimately result in phylogenomic datasets that can strongly mislead. For example, incongruence in the phylogeny of modern birds developed by Jarvis *et al.* (2014; 48 whole genomes) and Prum *et al.* (2015; 259 anchored phylogenomics loci, 198 species) has recently been attributed to differences in marker type rather than number of taxa (Reddy *et al.*, 2017). Whereas Jarvis *et al.* (2014) used primarily noncoding loci because they observed gross incongruence when using coding regions, the loci used by Prum *et al.* (2015), although nominally focused on broadly “anchored” conserved regions, came primarily from coding regions. Thus, at least one marker type is likely inappropriate when applied across modern birds (ca. 100 Ma). These data type effects can stem from multiple sources. Selection on exons might lead to localized differences in effective population size across the genome, as previous studies have highlighted issues with base composition heterogeneity within exons across taxa (Figuet *et al.*, 2015; Scally *et al.*, 2012). On the other hand, alignment quality of introns and ultraconserved elements can sometimes be less than desired (Edwards, Cloutier & Baker, 2017). Clearly marker effects can potentially have substantial consequences on species tree estimates and need to be further evaluated and compared side-by-side by the phylogenetic community (Shen, Hittinger & Rokas, 2017).

A priori versus a posteriori selection of loci for phylogenomics

In an ideal world, phylogeneticists would have whole and fully annotated genomes of all taxa available, allowing them to select loci for phylogenomics based on the relative merits of different loci. This a posteriori (i.e., after data generation) selection would be

carried out after assessing desired phylogenetic and biological properties of a wide array of markers for which the data are already in hand. A posteriori selection of loci for phylogenomics is clearly a long-term goal that will yield greater choice and justification for specific loci. Today, the loci for phylogenomics are selected a priori (i.e., before data generation) based primarily on cost and ease of collection and alignment, disregarding potentially useful regions of the genome. Thus, an attractive aspect of whole-genome sequencing (WGS) for phylogenomics is to have the opportunity to select markers a posteriori once genomes are in hand (e.g., [Edwards, Cloutier & Baker, 2017](#); [Fig. 1](#)). WGS is less prone to sequencing biases and also allows for further expansion into different research fields and questions based on the same initial data ([Lelieveld et al., 2015](#)). In contrast, a priori marker selection often limits the kinds of questions and methods that researchers can apply and represents a real constraint for phylogenomics and other disciplines.

An important constraint for using WGS for downstream phylogenomic analyses is genome quality. Obtaining high-coverage well-assembled and thoroughly annotated genomes is still very expensive and time-consuming, and even low-coverage genomes are still outside reach for large portions of the scientific community. However, even low-coverage genomes, which should be cautiously used for marker selection due to potential problems caused by poor annotation and coverage, can sometimes yield a modest number of markers for phylogenomics, and in the short term might even yield data sets allowing a broader diversity of markers for analysis. Although we are fully aware of these constraints, we are particularly excited about the potential that we see in routinely using WGS to produce phylogenomic data sets.

More taxa versus more loci

The question of whether to add more genes or more taxa was a dominant theme in phylogenetics in the 1990s and early 2000s (e.g., [Hillis, 1996](#); [Kim, 1996](#)), and remains a persistent question in phylogenomics today. After much debate in the literature (e.g., [Hillis, 1996](#); [Graybeal, 1998](#); [Hillis, 1998](#); [Poe, 1998](#); [Mitchell, Mitter & Regier, 2000](#)), the initial consensus view from the Sanger sequencing era of phylogenetics, is that adding more taxa generally improves phylogenetic analysis more so than more markers (e.g., [Hillis, 1996](#); [Graybeal, 1998](#); [Poe, 1998](#)). However, phylogenomics is adding a new twist to this consensus, both from the standpoint of data acquisition and from theory (e.g., [Rokas & Carroll, 2005](#); [Nabhan & Sarkar, 2012](#); [Xi et al., 2012](#); [Patel, Kimball & Braun, 2013](#)). Amassing large data sets, both in terms of more taxa and more loci, is still a guiding principle of phylogenomics. But with the ability now to bring together many different types of markers in a single analysis, and to analyze them in ways that were not previously available, the “more taxa vs. more genes” debate is becoming more nuanced ([Nabhan & Sarkar, 2012](#)). For example, recent work shows that this debate can be highly context-specific and model-dependent (e.g., [Baurain, Brinkmann & Philippe, 2006](#); [Dell Ampio et al., 2013](#); [Edwards et al., 2016b](#)). Also, coalescent methods appear to be more robust to limited taxon sampling than traditional methods like concatenation ([Song et al., 2012](#); [Liu, Xi & Davis, 2015](#)). Some researchers favor “horizontal” data

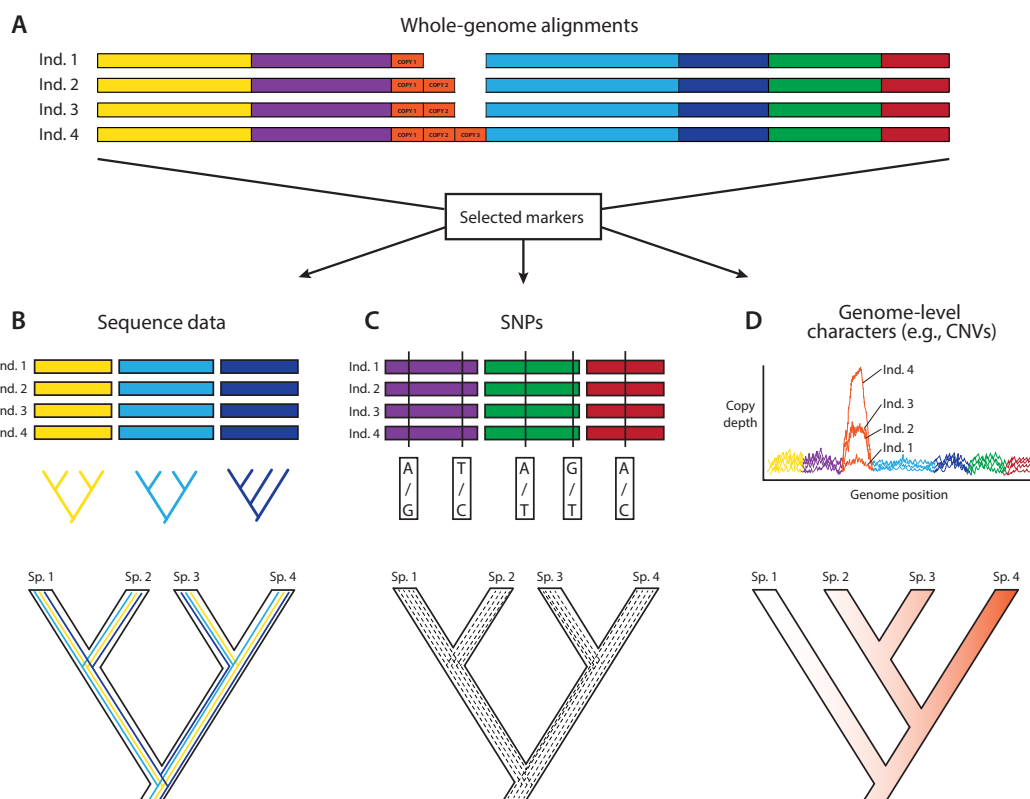


Figure 1 A posteriori marker selection from whole-genome alignments for phylogenomics and phylogeography. Whole-genome analysis (A) permits researchers to choose different markers for specific purposes (B–D). By contrast, subsampling methods such as Rad-seq or hybrid capture, which dominate phylogenomics today, usually yield a specific set of markers that the researcher has chosen a priori. The generation of WGA thus greatly increases the use of genomic data in biological research, beyond the initial goals of the researcher producing those data. Here, we show how a hypothetical WGA that includes seven different loci (different colors) for four individuals allows extracting sequence data to generate gene trees (B), identifying SNPs to genotype individuals (C), and measuring copy depth to infer CNVs across genomic regions (D). Ultimately, these different kinds of data can be translated into species tree inferences (B–D). In the case of CNVs, only locus number 3 (orange) shows significant CNV. Because CNVs are measured as continuous characters (i.e., copy depth), the orange shading represents a hypothetical evolutionary scenario of copy number variation of genomic region number 3 within the inferred species tree, which is incongruent with those based on sequence and SNP data from other loci in the genome.

Full-size DOI: 10.7717/peerj.6399/fig-1

matrices, wherein the number of loci far exceeds the number of taxa, whereas other researchers favor “vertical” matrices, where many taxa are analyzed at just a few (1–5) loci. Whereas the PCR era of phylogenetics was often dominated by vertical matrices, HTS is allowing data matrices to become more horizontal (Fig. 2). It is important to note that as these more horizontal data become more prevalent, they increase the amount of missing data and aligning problems that can contribute to misleading or low phylogenetic resolution. At least in a coalescent framework, scaling up in both dimensions will be crucial for improved phylogenies and the number of loci required to resolve a given phylogenetic problem is often a function of the coalescent branch lengths in the phylogenetic tree being resolved, with longer branches requiring fewer loci (Edwards, Liu & Pearl, 2007; Huang et al., 2010). At deeper time scales, increasing the number of loci have not proven

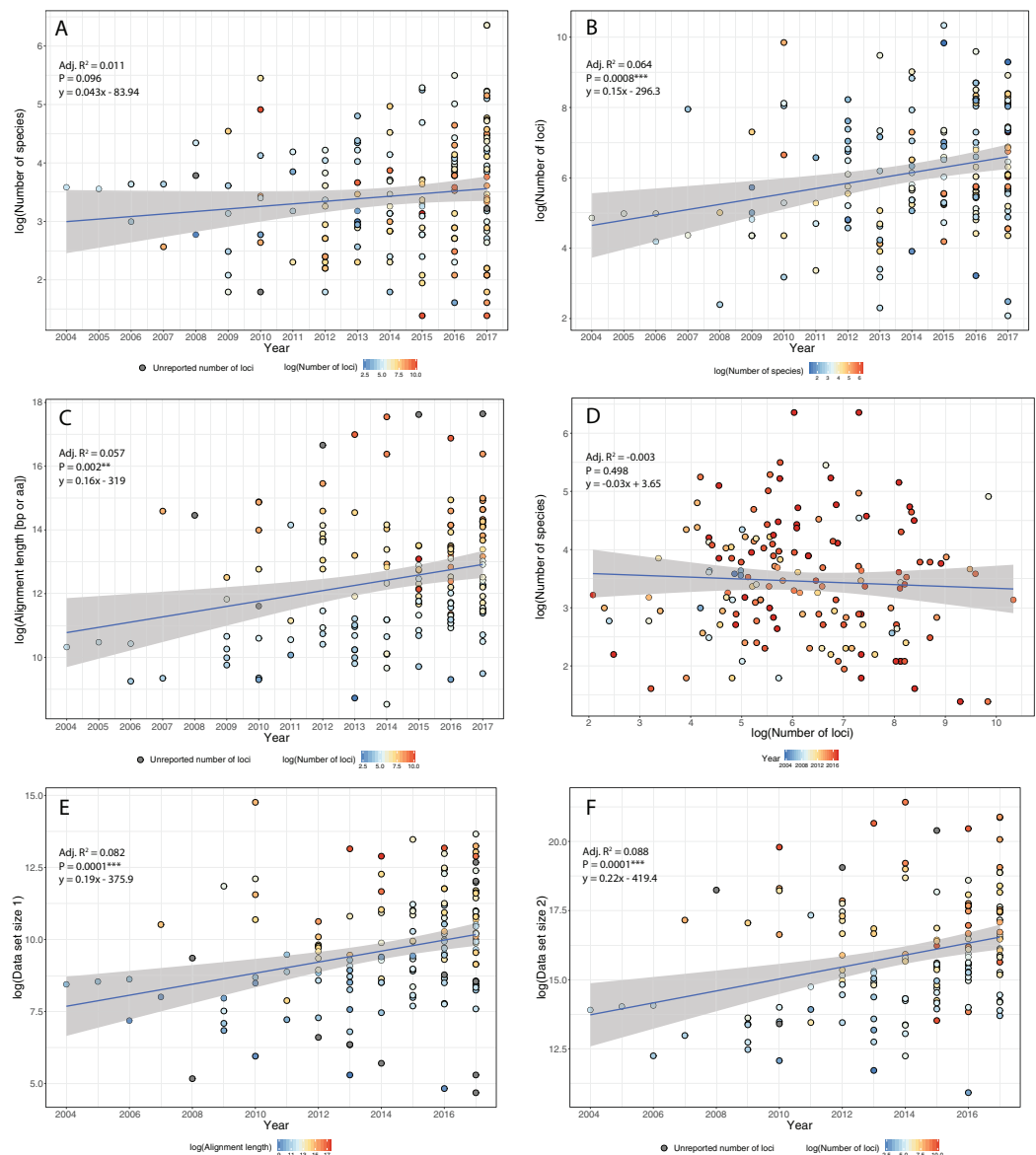


Figure 2 Trends in phylogenomic data sets since the emergence of HTS. Based on a sample of 164 phylogenomic papers published since 2004 (see Table S1), we observed no increase in the number of species per data set over time (A). On the other hand, there is a significant increase in the number of loci (B), total alignment length (C), and total data set size, as measured by the product of species times locus number (Data set size 1, E) and species times total alignment length (Data set size 2, F). Moreover, the advent of HTS does not support the notion of a tradeoff between the number of species and the number of loci in phylogenomic studies (D).

Full-size [DOI: 10.7717/peerj.6399/fig-2](https://doi.org/10.7717/peerj.6399/fig-2)

particularly useful to resolve problematic areas in the ToL (e.g., sister group to animals; [King & Rokas, 2017](#)). Despite methodological complications in the accurate estimation of population genetic parameters and computational limitations for MCMC convergence at deep times in coalescent-based analyses, recalcitrant nodes likely represent true complexities in the diversification history of these groups and not necessarily reflect failures of coalescent-based phylogenetics ([Lanier & Knowles, 2015](#)).

To study how researchers have resolved challenges of balancing numbers of taxa versus numbers of loci, we quantified trends in phylogenomic data set size and structure over the past 13 years, drawing data from 164 data sets across diverse taxa (Table S1). We found that, whereas the number of species per paper has not increased significantly over time (Fig. 2A), there were significant increases with time in number of loci (Fig. 2B), total length of sequence analyzed (Fig. 2C), as well as total data set size, as measured by the product of species times locus number (Fig. 2E) or species times total alignment length (Fig. 2F). These results mirror similar trends evaluated for the size of data sets in phylogeography (Garrick *et al.*, 2015). Surprisingly, we found no evidence for a tradeoff between the number of species investigated and the number of loci analyzed (Fig. 2D); perhaps HTS data sets have plateaued somewhat in terms of number of loci, whereas the number of species analyzed is more a function of the questions being asked and the clade being investigated. Regardless, we suspect that, in general, the number of loci and total alignment lengths in phylogenomic data sets are likely a function of resources and sequencing effort. The era of whole-genome sequencing in phylogenomics is still dawning, given that most studies thus far have used targeted approaches for sampling loci (Table S1). We suspect that once whole-genome sequencing on a clade-wide basis become routine (e.g., *Genome 10K Community of Scientists*, 2009; Grigoriev *et al.*, 2014; Cheng *et al.*, 2018), we will witness yet another jump in the sizes of phylogenomic data sets.

Filtering heterogeneous phylogenomic data sets

Recent studies show that the addition of more loci and more taxa can result in higher levels of gene-tree discordance (e.g., Smith *et al.*, 2015; Shen, Hittinger & Rokas, 2017). This is not unexpected - as the number of taxa and loci increase, the greater the likelihood of capturing signals of the heterogeneous evolutionary history (e.g., ILS, lateral gene transfer (LGT), hybridization, gene duplication and loss (GDL); See Table 1 for a definition of these concepts), misidentifying orthologs from paralogs, and recovering patterns of molecular evolution (e.g., noise/lack of signal in the sequences, and nonstationarity in base composition) that can contribute to gene tree discord. At the same time, the variance in gene tree topologies could also have been caused by errors in gene tree estimation. Such observations have been used to argue that the accuracy of gene tree inference should be maximized or at least evaluated, but it is not clear what criteria should be used to filter sets of gene trees. For example, filters can be based on rates of molecular evolution (Klopfstein, Massingham & Goldman, 2017), levels of phylogenetic informativeness (Fong *et al.*, 2012), or on the cause of gene-tree discord itself, if known (Huang *et al.*, 2010). Chen, Liang & Zhang (2015) found that selecting genes whose trees contained a well-known uncontested and long branch in a given species phylogeny (long enough so as not to incur substantial ILS) was a better way to improve phylogenomic signal than selecting genes based on characteristics of sequence evolution. All of these methods are excellent suggestions and should be explored further. Still, the effects of such culling on the distribution of gene trees, and whether it could distort the distribution so that it departs from models like the multispecies coalescent, are unknown, and potentially of

Table 1 Definitions of core concepts used in this article.

| Concept | Definition |
|---|--|
| The Tree of Life (ToL) | This idea, originally articulated by Darwin and others, refers to the grand vision of understanding the branching pattern of all life on earth. Today the idea conveys the use of morphological and molecular data to reconstruct the phylogenetic relationships of all life forms. In some usages, the idea also includes reconstructing reticulate evolutionary events, such as introgression and hybridization, which are now thought to be common in many lineages. |
| High-throughput sequencing (HTS) | Also referred to as “next generation sequencing”, this term refers to the plethora of new DNA and RNA sequencing technologies that in the last fifteen years have allowed biologists to dramatically increase the number of bases sequenced for a given species or clade. HTS technologies can be applied to sequencing whole genomes or transcriptomes and have been embraced by phylogeneticists interested in increasing the size of comparative molecular data sets. See Goodwin, McPherson & McCombie (2016) for a review on the progress of HTS. |
| The multispecies coalescent model (MSC) | A generalization of the standard, single population coalescent model to multiple species related in a phylogeny. The MSC applies the single-population coalescent model to each branch of a phylogenetic tree, including both terminal and internal branches. In the MSC, alleles sampled in terminal species will coalesce to a smaller number of ancestral alleles at a rate depending on the effective population size within the branch. The gene tree lineages in a branch of the species tree do not necessarily coalesce within that branch as one goes backwards in time; multiple alleles may persist into ancestral branches. This phenomenon is called incomplete lineage sorting (see next definition). The decrease in the number of alleles and the time to coalescence to a single allele in a lineage follows the standard neutral coalescent model, until all alleles coalesce from all species. See Rannala & Yang (2003) and Degnan & Rosenberg (2009) for a full discussion. |
| Incomplete lineage sorting (ILS) | This phenomenon, originally described by John Avise (see Avise et al. 1987) refers to the tendency of alleles in an ancestral species to persist across multiple speciation events, resulting in a situation in which the gene tree (“allele tree”) differs from the species tree. In ILS, alleles fail to “sort” by genetic drift as species diverge from one another, resulting in different species retaining the same alleles, or their descendants, causing discordance with the overarching species or population tree. The language of this phrase is looking forward in time, as opposed to the language of coalescence, which looks backwards in time. See Degnan & Rosenberg (2009) for a full discussion. |
| Gene duplication and loss (GDL) | This concept describes the process by which a gene in an ancestral species can duplicate, forming paralogs and one or more of the paralogs can subsequently be deleted from the genome, resulting in complex patterns of relationships among paralogs and orthologs. Gene duplication is another mechanism, in addition to ILS, that can render the gene tree different from the species tree. As a result of gene (paralog) loss, inferring the correct ortholog/paralog relationships and history of branching events in a multigene family can be challenging. Phylogenetic models incorporating GDL try to use patterns in multigene families to deduce the branching history of the constituent species. See Degnan & Rosenberg (2009) and Sousa et al. (2017) for a full discussion. |
| Ancestral recombination graph (ARG) | This is a complete record of the coalescent and recombination events in the history of a set of DNA sequences. As a consequence of incorporating recombination events, ARGs do not necessarily depict trees, but often have a network structure. The accurate estimation of ARGs remains challenging but they enhance our ability to estimate recombination rates, ancestral effective population sizes, population divergence times, rates of gene flow between populations, and detect selective sweeps. See Griffiths & Marjoram (1996) , Siepel (2009) , and Rasmussen et al. (2014) for a full discussion. |
| Lateral gene transfer (LGT) | This process occurs when genes jump taxonomic and phylogenetic boundaries, moving between unrelated species and therefore causing discordances between genetic and lineage history. LGT, along with ILS and GDL was among the three primary causes of discordance between gene and species trees identified by Maddison (1997) . LGT has been documented to occur between bacterial lineages, between bacteria and viruses, and between these two and eukaryotes, including plants and vertebrates. If not identified prior to phylogenetic analysis, LGT can cause many algorithms for phylogenetic inference to fail. Without prior identification, LGT essentially amounts to errors in data sets and sequence alignments. At the same time, LGT can be a source of adaptation and evolutionary novelty for recipient genomes and has had a major impact on the history of life. See Gogarten & Townsend (2005) for a full review. |

concern (but see [Huang et al., 2017](#)). Aside from the use of some explicit methods for detection of outlier genes (e.g., [de Vienne, Ollier & Aguileta, 2012](#)), rogue taxa (e.g., [Aberer, Krompas & Stamatakis, 2013](#)), outlier long branches (e.g., [Mai & Mirarab, 2018](#)), and tree space visualization (e.g., [Huang et al., 2016](#); [Jombart et al., 2017](#)), an obvious way to alleviate potential effects of gene tree outliers is a more balanced taxon sampling

(Hedtke, Townsend & Hillis, 2006). Nonetheless, we need further studies on the effects of different types of phylogenomic filters on the properties of large-scale phylogenomic datasets.

High-throughput sequencing opens possibilities for new information and marker types

Heterozygosity and intra-individual site polymorphisms

Some of the prevalent occurrences in organisms with multiple ploidies are intra-individual polymorphisms and increased heterozygosity. However, due to issues such as lack of sufficient read coverage and connectivity, confident identification of such polymorphism continues to be challenging (Garrick, Sunnucks & Dyer, 2010; Lischer, Excoffier & Heckel, 2014; Schrempf et al., 2016) and many data sets do not permit statistical approaches, such as PHASE (Stephens, Smith & Donnelly, 2001), to robustly determine haplotypes of different alleles (Garrick, Sunnucks & Dyer, 2010). Consequently, in phylogenetics, heterozygosity and intra-specific polymorphic sites are often accommodated using UIPAC ambiguity codes or ignored entirely or by randomly selecting alleles (Iqbal et al., 2012). In fact, most “one sequence per individual/species” phylogenomic data sets consists of haplotypes that might not occur in nature because many methods, including de novo assemblies of genomes, yield single haplotypes consisting of consensus or other haplotype summaries from diploid organisms. The fact that HTS produces several reads of the same region allows for the identification of heterozygosity and intra-specific polymorphic sites represents an untapped opportunity to incorporate intra-individual variation in our phylogenetic estimates (Lischer, Excoffier & Heckel, 2014; Schrempf et al., 2016; Andermann et al., 2018). Recent models have been proposed to improve calling and sorting such polymorphisms (De Maio, Schlötterer & Kosiol, 2013; Lischer, Excoffier & Heckel, 2014; Potts, Hedderson & Grimm, 2014; Schrempf et al., 2016) and, although results of different studies vary (Kubatko, Gibbs & Bloomquist, 2011; Lischer, Excoffier & Heckel, 2014), estimation of individual, naturally occurring haplotypes has been shown to improve phylogenomic reconstructions based on genome-scale data (Andermann et al., 2018).

Rare genomic changes

As noted above, molecular phylogenetics has primarily used alignments of sequence-level data for phylogenetic inference. This bias is perhaps driven by the notion that genome evolution occurs by aggregating small changes, such as point substitutions, over time. However, this bias also responds to the challenges of characterizing rare genomic changes, such as indels, transpositions, inversions, and other large-scale genomic events (Rokas & Holland, 2000; Boore, 2006; Bleidorn, 2017). This emphasis on sequence data has produced a vast ecosystem of algorithms tailored to analyze such data, but most phylogeneticists would agree that rare genomic changes would be a welcome addition to the toolkit of phylogenomics, because they are generally regarded as highly informative markers, providing strong evidence of homology and monophyly (Boore, 2006; Rogozin et al., 2008). With the increased availability and affordability of WGS, our view of genome plasticity has changed drastically in recent years and we are now capable of

exploring other genomic features beyond the signals encapsulated in DNA or amino acid sequences (e.g., [Ryan et al., 2013](#)). The question then arises of how to identify and utilize these rare genomic markers, as well and assess their phylogenetic informativeness ([Rokas & Holland, 2000](#); [King & Rokas, 2017](#)). Genome-level characters will likely have different evolutionary properties than sequence-based markers, suggesting that one of the biggest challenges we face for incorporating genomic changes into phylogenetic analyses is to find informative evolutionary models and tools suited for these kinds of data and assess how congruent or discordant they are with respect to other markers (e.g., [Rota-Stabelli et al., 2011](#)). This will not only shed light on how phylogenetically informative different genomic changes are, but also will broaden our understanding of the evolutionary intricacies across different genomic regions ([Rokas, 2011](#); [Leigh et al., 2011](#)).

Gene order and synteny

Computational algorithms to use gene order and rearrangements as markers in phylogenetics ([Tang et al., 2004](#); [Ghiurcuta & Moret, 2014](#); [Kowada et al., 2016](#)) were spurred in part by the seminal paper by [Boore, Daehler & Brown \(1999\)](#) using mitochondrial gene rearrangements to understand the phylogeny of arthropods. Initially, algorithms for making use of gene order and synteny were applied primarily to microbial genomes, but recent efforts have extended such methods to the analysis of eukaryotes as well (see [Lin et al., 2013](#)). Gene order and synteny appear most promising at high phylogenetic levels, although we still do not know how informative gene order will be at many levels. For instance, chromosomal rearrangements appear highly dynamic in some groups, such as mammals, and further study of their use in phylogenomics is warranted ([Murphy et al., 2005](#)).

Indels and transpositions

Indels and transpositions are two types of molecular characters that are underutilized in phylogenomics. The former perhaps because standard methods of analysis often treat indels as missing data and the latter because they are technically challenging to collect without whole-genome data. Indels have been used sporadically in phylogenomics and several researchers have argued for their utility and informativeness, given appropriate analytical tools ([Jarvis et al., 2014](#); [Ashkenazy et al., 2014](#); [Roncal et al., 2016](#)). [Murphy et al. \(2007\)](#) used indels in protein-coding regions to bolster estimates of mammalian phylogeny and found that the Atlantogenata hypothesis was supported after scrutinizing proteome-wide indels for spurious alignments and orthology. The Avian Phylogenomics Project ([Zhang, Jarvis & Gilbert, 2014](#)) found that indels had less homoplasy than SNPs and, despite showing high levels of ILS, was largely congruent with other markers across the avian tree. Transposable elements arguably are even more highly favored by phylogenomics researchers, but are much more difficult to isolate and analyze and have been used principally across various studies in mammals and birds ([Kaiser, van Tuinen & Ellegren, 2007](#); [Churakov et al., 2010](#); [Kriegs et al., 2010](#); [Suh et al., 2011](#); [Baker et al., 2014](#); [Cloutier et al., 2018a](#)). Whereas they are generally considered to have a low rate of homoplasy, most researchers agree that they can in some

circumstances exhibit insertional homoplasy. Moreover, no marker is immune to the challenges of ILS, and transposable elements and indels are no exception ([Matzke et al., 2012](#); [Suh, Smeds & Ellegren, 2015](#)). Still, the exceptional resolution afforded by some studies employing transposable elements is exciting, and we expect this marker type to increase in use as whole genomes are collected with higher frequency.

Copy number variations

The 1000 Genomes Project estimates that in humans about 20 million base pairs are affected by structural variants, including copy number variations (CNV) and large deletions ([1000 Genomes Project Consortium, 2015](#)), suggesting that these types of mutations encompass a higher fraction of the human genome than do SNPs. A CNV is a DNA segment of at least one kilobase (kb) that varies in copy number compared with a reference genome ([Redon et al., 2006](#)). CNVs appear as deletions, insertions, duplications, and complex multi-site variants ([Fredman et al., 2004](#)). Such a profusion of CNVs across human genomes has proven useful in tracking population structure ([Sjödén & Jakobsson, 2012](#)), but still remains underappreciated in phylogenetics.

Newly available methods allow inference of CNV at high resolution with great accuracy ([Wiedenhoeft, Brugel & Schliep, 2016](#)). The frequency with which CNVs occur in animal and plant populations raises the question of how informative they would be at higher phylogenetic levels, and whether they would incur unwanted homoplasy that would obscure homology and phylogenetic relationships. For example, some CNVs evolve so quickly that they can be used with success at the sub-individual level, for example, in tracking clonal evolution of cancer cells using CNV-specific phylogenetic methods (e.g., [Schwartz & Schäffer, 2017](#); [Liang, Liao & Zhu, 2017](#); [Ricketts et al., 2018](#); [Urrutia et al., 2018](#)). Moreover, their interspecific variation has been shown to correlate with the phylogeny of some groups, such as the highly pathogenic and rapidly-evolving barley powdery mildews ([Frantzeskakis et al., 2018](#)). However, such fast evolution may mean that these markers might be less useful at higher levels of biological organization. Additionally, the adaptive nature of CNVs may or may not facilitate clear phylogenetic signals. For example, a study in *Arabidopsis thaliana* showed that adaptation to novel environments, or to varying temperatures, is associated with mutations in CNVs ([DeBolt, 2010](#)). If CNVs are to become a useful tool in phylogenomics or phylogeography, we must understand their microevolutionary properties in greater detail. For example, the pattern of evolution of CNVs, wherein deletions of genetic material may not easily revert, resulting in a type of Dollo evolution, might help clarify the overall structure of the models applied to them ([Rogozin et al., 2006](#); [Gusfield, 2015](#)).

Recent advances in the generation of high-throughput sequence data and their impact on the reconstruction of the Tree of Life

As sequencing technology rapidly moves forward (reviewed by [Goodwin, McPherson & McCombie, 2016](#)), our ability to accurately identify the aforementioned marker types increases considerably. For instance, the low per-base error rate of short-read sequencing technologies, such as the Illumina HiSeq X Ten and NovaSeq (Illumina, San Diego,

CA, USA), allow for a significant reduction in the cost of sequencing which can result in data for more taxa at a higher coverage. This is certainly beneficial for the accurate identification of SNPs, heterozygosity, and intra-individual site polymorphisms (Goodwin, McPherson & McCombie, 2016) and their use in a phylogenomic context. Moreover, single molecule real-time sequencing technologies, such as Pacific Biosciences (Pacific Biosciences of California, Menlo Park, CA, USA) and Oxford Nanopore (Oxford Nanopore Technologies, Oxford, UK) produce reads that exceed 10 kb (Rhoads & Au, 2015; Lu, Giordano & Ning, 2016). The advent of these technologies has led to improved and more efficient assembly methods that allow accurate identification of structural changes (e.g., Khost, Eickbush & Larracuenta, 2017; Merker et al., 2018). When combined with data resulting from short-read sequencing, they represent a powerful tool to correctly identify and use a wide array of genomic markers for numerous purposes. These technical advances, which include portable devices that can be carried into the field (i.e., Oxford Nanopore; Johnson et al., 2017), will certainly yield an increase in the genomic loci and taxa available for genomic and phylogenomic studies.

CONCEPTS AND MODELS IN PHYLOGENOMICS

For decades, phylogenetics has struggled with how best to translate evolutionary changes in DNA sequences and other characters into phylogenies, and genomic data are no exception to this trend. Phylogenomics is still in a developing stage of formulating models that effectively represent the underlying mechanisms for genome-scale variation while remaining efficient and within reasonable analytical and bioinformatic capacities. The current focus on models and evolutionary forces generating the patterns that we recover as branching and reticulation events in our phylogenetic reconstructions is a healthy one and can be extended to other important topics in phylogenomics, such as species delimitation, character mapping, and trait evolution (e.g., Yang & Rannala, 2014). All of these areas are developing rapidly and are in need of updated models and bioinformatics applications to cope with the heterogeneity brought by genome-scale data.

The multispecies coalescent model

One of the key practical advances in molecular phylogenetics has been the incorporation of gene tree stochasticity into the inference of species phylogenies, via the multispecies coalescent model (MSC: Rannala & Yang, 2003; Liu & Pearl, 2007; Heled & Drummond, 2010). The MSC allows gene trees to be inferred with their own histories, including coalescent-appropriate branching models, but contained within independent yet connected lineages within a species phylogeny, with speciation-appropriate branching models (Degnan & Rosenberg, 2009). The main conceptual advance has been to understand and separately manage the variation at different levels of biological organization—an advance that began years ago (Doyle, 1992; Maddison, 1997; Pamilo & Nei, 1988), but has only recently been widely embraced and put into practice (Edwards, 2009b). Given its ability to accommodate heterogeneous histories across loci scattered throughout the genome, the MSC lays at the core of the conceptual framework

to deal with genome-scale data (e.g., [Rannala & Yang, 2008](#); [Liu et al., 2015](#)). In the few instances in which model comparison and fit has been evaluated ([Liu & Pearl, 2007](#); [Edwards, Liu & Pearl, 2007](#)), the MSC vastly outperforms concatenation. This of course does not mean that the MSC is the correct, or even an adequate, model for phylogenomic data ([Reid et al., 2014](#)). Despite concerns regarding some of its implementations when dealing with genomic data (e.g., [Springer & Gatesy, 2016](#)), the MSC is a powerful theoretical model for phylogenomics and there is room for refinement and improvement for its applications (e.g., [Edwards et al., 2016b](#), [Xu & Yang, 2016](#)).

Bypassing full likelihood models by relying on summaries of the coalescent process

Given the huge computational resources required for modelling all the complexities of evolutionary processes in a statistical framework, there is interest in methods that will accommodate genome-scale data for large numbers of species within a coalescent framework. The utility of such methods cannot be overstated: the rapid rise of large-scale genomic data sets has clearly outstripped theoretical and computational methods required to analyze them. For example, although progress is being made regarding scalability of full Bayesian methods of species phylogeny inference (e.g., [Ogilvie, Bouckaert & Drummond, 2017](#)), they are still unable to accommodate large phylogenomic datasets, which often consist of hundreds of species for thousands of loci ([Table S1](#)). A common approach to speeding up species phylogeny inference consists of ‘two-step’ methods, wherein gene trees are estimated first and separately from the species phylogeny; then, using various summaries of the coalescent process for collections of gene trees, a species phylogeny is estimated. Many useful methods for estimating species phylogenies in this way have been proposed (see [Marcussen et al., 2014](#); [Liu, Wu & Yu, 2015](#); [Mirarab & Warnow, 2015](#); [Mirarab, Bayzid & Warnow, 2016](#)), taking advantage of summary statistics of the coalescent process, such as the average ranks of pairs of species in the collection of gene trees (e.g., STAR: [Liu et al., 2009](#); ASTRAL-II: [Mirarab & Warnow, 2015](#)) or the distribution of gene trees containing triplets of species (e.g., MP-EST; [Liu, Yu & Edwards, 2010](#)). Some of these two-step methods, while approximate, nonetheless allow for statistical testing in a likelihood framework. For example, MP-EST can evaluate the (pseudo)likelihood of two proposed species phylogenies given a collection of gene trees and the difference in likelihood can be used to evaluate two proposed species phylogenies against each other. However, such statistical approaches have rarely been used thus far, and bootstrapping or approximate posterior probabilities on branches are by far the most common statistics applied to species phylogenies ([Sayyari & Mirarab, 2016](#)). Speeding up the estimation process using two-step methods can be effective, but it can also accumulate errors or misallocate sources of variance which cannot be corrected at later stages ([Xu & Yang, 2016](#)). If gene trees are biased or uninformative, then downstream analyses for species phylogeny estimation or species delimitation may similarly be compromised (e.g., [Olave, Sola & Knowles, 2014](#)). For example, MP-EST can sometimes perform poorly when PhyML ([Guindon et al., 2010](#)) is used to build low-information gene trees because PhyML may produce biased gene trees when the alignments contain

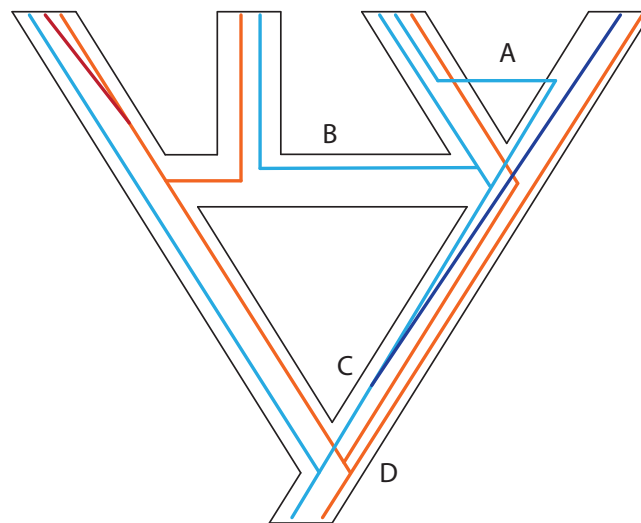


Figure 3 Some examples of violations of the multispecies coalescent. In event A, there is gene flow; in event B there is homoploid hybridization; in event C, there is a gene duplication; and in event D, incomplete lineage sorting. All of these processes contribute to gene tree heterogeneity but fall outside the standard multispecies coalescent model. Importantly, all of these processes also yield strictly dichotomous gene trees, whereas recombination (not illustrated here) does not.

Full-size DOI: 10.7717/peerj.6399/fig-3

very similar sequences (Xi, Liu & Davis, 2015). This may account for the lower performance of MP-EST compared to ASTRAL in some simulation conditions, because ASTRAL resolves input polytomies and zero-length branches in gene trees more appropriately. This difference between MP-EST and ASTRAL is eliminated when RAxML (Stamatakis, 2014) is used to build gene trees (Xi, Liu & Davis, 2015).

Beyond the multispecies coalescent model

Reticulation at multiple levels challenges the standard multispecies coalescent model

The phylogenetic processes of branching and reticulation can operate at several levels of organization, including within genes, within genomes, and within populations or species (Figs. 3–5). For example, recombination can cause reticulations within genes, allopolyploidization can cause reticulations at the level of whole genomes, and introgression and hybridization can cause reticulations at the level of populations. These levels are nested so that branching processes (and in part reticulations) acting at a higher level will cause correlated branching patterns at lower levels. At the same time, reticulations at lower levels, such as recombination acting within genes, will cause inference problems at higher levels, such as estimating population histories. Crucially, however, it is only recombination that will break one key element driving many recent models of phylogenetics and population histories, namely dichotomous gene trees. Reticulations at levels of organization higher than the genome, such as the fusing of populations, as well as gene duplication, will still yield collections of dichotomous gene trees, even if the higher-level history is reticulated. Ultimately, the additive effects of

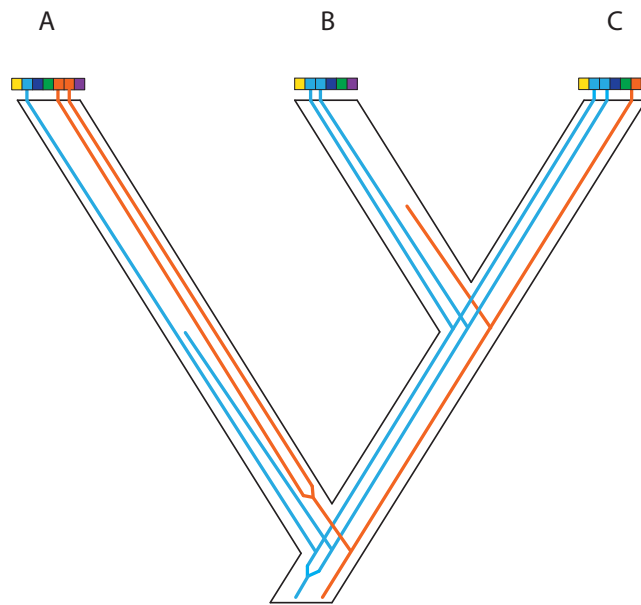


Figure 4 Gene duplication and loss (GDL) creates patterns that can mimic incomplete lineage sorting and other processes, leading to spurious inferences of the species history. Genes and genomes of three species A, B, and C. Multi-colored bars show (parts of) their genomes with a number of loci indicated in different colors. The orange gene is duplicated in species A and it was lost in species B. The blue gene was duplicated before the divergence between species A and the ancestor of species B and C. However, one of these copies was lost in species A, whereas both copies were maintained in species B and C. Reconstruction of the orange gene tree based on extant diversity will yield a wrong inference of its history due to the absence of data for species B. On the other hand, a phylogenetic reconstruction of the blue gene is difficult to predict. Depending on which of the duplicates are sampled for species B and C, different outcomes can be expected regarding the relationship among the three species. The duplication and loss history of these two genes may cause serious issues for phylogenetic reconstruction because no specific pattern can be expected between them.

Full-size [DOI: 10.7717/peerj.6399/fig-4](https://doi.org/10.7717/peerj.6399/fig-4)

these reticulate processes result in our observed phylogenetic reconstructions, and we expect all of these scenarios to produce bifurcating, dichotomous gene trees. From a modelling point of view, another key distinction is whether at the species level, we still have a phylogeny that is tree-like, or whether a network is needed. The process whereby two populations jointly produce a third requires a network to model properly. Allopolyploidy is another situation requiring a network. There are several statistical methods for inferring homoploid networks (Yu *et al.*, 2014; Solís-Lemus & Ané, 2016; Wen, Yu & Nakhleh, 2016; Wen & Nakhleh, 2018), species histories under allopolyploidy (Jones, Sagitov & Oxelman, 2013), and some two-step methods such as PADRE (Huber *et al.*, 2006; Lott *et al.*, 2009). In general, dealing with multiple simultaneous violations of the MSC, such as introgression and allopolyploidy, remains challenging (Degnan, 2018). It is likely that the history of many radiations involves parts of the genome with a dichotomous history and parts that exhibit reticulation, demanding methods that accommodate both scenarios. Alternatively, rather than trying to accommodate multiple processes in our methods for phylogenetic inference, we might instead focus our attention on subsets of loci that would not violate the MSC (e.g., Knowles, Smith & Sukumaran, 2018). In cases where processes other than ILS contribute to gene tree

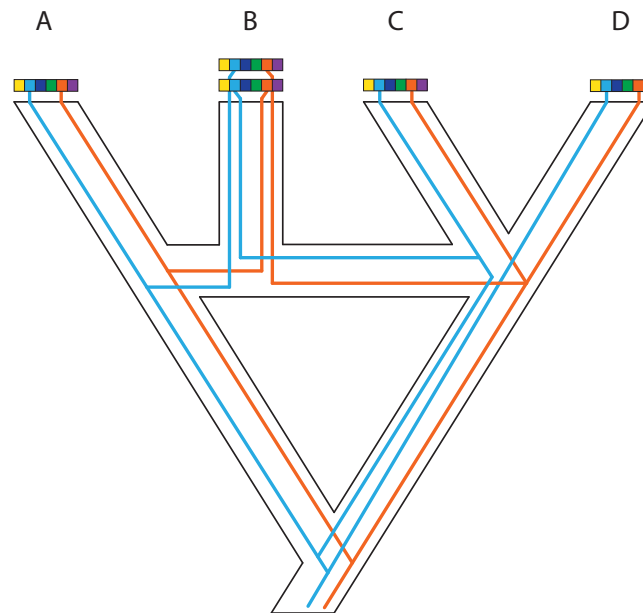


Figure 5 Complex patterns of gene lineages with polyploidization and interspecific gene flow. Genes and genomes of four species A, B, C and D. Multi-colored bars show (parts of) genomes with a number of loci indicated in different colors. Two gene trees, one orange and one blue, evolve within the species network. Species B is an allopolyploid containing two genomes.

Full-size DOI: 10.7717/peerj.6399/fig-5

discord (i.e., the distribution of trees is statistically inconsistent with expectations under the MSC; see [Smith et al., 2015](#)), loci consistent with the MSC can be identified (e.g., separated from loci with horizontal gene transfer) using CLASSIPHY ([Huang et al., 2017](#)). It has also been suggested that in order to distinguish violations of the MSC, the MSC can be used as a null model to be compared with increasingly complex models that would invoke processes such as hybridization and recombination using networks ([Degnan, 2018](#)). To follow this promising approach, further research must be conducted to not only model specific processes, but also distinguish them.

Models accommodating dichotomous divergence with gene flow are somewhat limited. For example, in IMA2 ([Hey & Nielsen, 2004](#); [Hey & Nielsen, 2007](#); [Hey, 2010](#)) the species phylogeny must be known and fairly small; in the method of [Dalquen, Zhu & Yang \(2017\)](#), both the species phylogeny and gene trees are restricted to three tips. Looking forward, it may be useful to deal with two sub-problems: The first is estimating the species phylogeny despite migration, for example by identifying which loci are interfering with the species phylogeny inference or causing reticulations in the form of gene flow. The second sub-problem is to incorporate a gradual speciation process ([Fig. 6](#)), where gene flow after speciation slowly declines, perhaps according to some simple function like an exponential. Such a model would capture what is thought to be a more common speciation process than the instantaneous process modelled by the MSC ([Jones, 2018](#)).

In some cases, it is possible to model one situation with either a species network or a tree with gene flow. [Long \(1991\)](#) discussed two models of admixture: Intermixture and gene flow ([Fig. 7](#)). The phylogenetics community has mainly focused on methods for

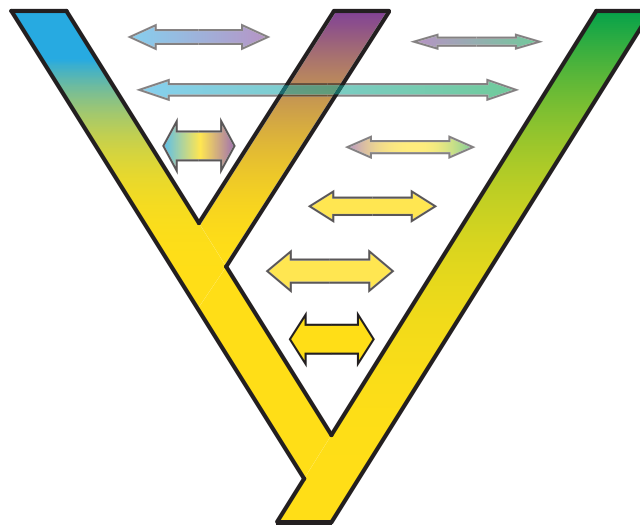


Figure 6 Gradual speciation, or isolation-with migration. After starting to split, gene flow between species decreases gradually. Such a gradual decrease in the extent of gene flow between species might present an especially useful extension of the standard multispecies coalescent model. Colors depict different gene pools and their gradual change along branches describes how species gradually differentiate despite the existence of migration over time. Thickness and color intensity of arrows show that gene flow becomes weaker as species gradually isolate. [Full-size](#) DOI: 10.7717/peerj.6399/fig-6

inference under the intermixture model (e.g., the multispecies network coalescent; [Yu et al., 2014](#), whereas the population genetics community has focused more on models including gene flow (e.g., IM ([Hey & Nielsen, 2007](#)), G-PhoCS ([Gronau et al., 2011](#)), PHRAPL ([Jackson et al., 2017](#)), admixture graphs)). While some initial work to test inference based on one of these models on data generated by the other has recently appeared ([Wen & Nakhleh, 2018](#); [Solís-Lemus, Bastide & Ané, 2017](#); [Blischak et al., 2018](#); [Zhang et al., 2018](#)), much more work is needed to bring together these two lines of work. Simulations and comparisons of observed and expected summary statistics, such as the site-frequency spectrum ([Excoffier et al., 2013](#)), have proven especially useful in distinguishing such scenarios (Fig. 7).

Reticulation in the form of gene flow or introgression is probably the most difficult violation of the MSC to address (but see [Hibbins & Hahn, 2018](#) for a model that estimates the timing and direction of introgression based on the multispecies network coalescent), in part because the number of potential trees accommodating a reticulating network is even higher than the already high number of trees for a given number of taxa. There is at least one issue where reticulation presents an opportunity as well as a challenge. Any kind of gene flow/hybridization means that there is the possibility of inferring the existence of extinct species, because extinct species contribute novel alleles that exceed the coalescence time of most alleles in the focal species under study ([Hammer et al., 2011](#)). Well-known examples are the documented presence of Neanderthal genes in most human genomes due to introgression (e.g., [Meyer et al., 2012](#)) and the presence of genomes derived from now-extinct diploids in extant allopolyploids (i.e. meso-allopolyploids; e.g., [Mandáková et al., 2010](#); [Marcussen et al., 2015](#)). Some current models can explain the

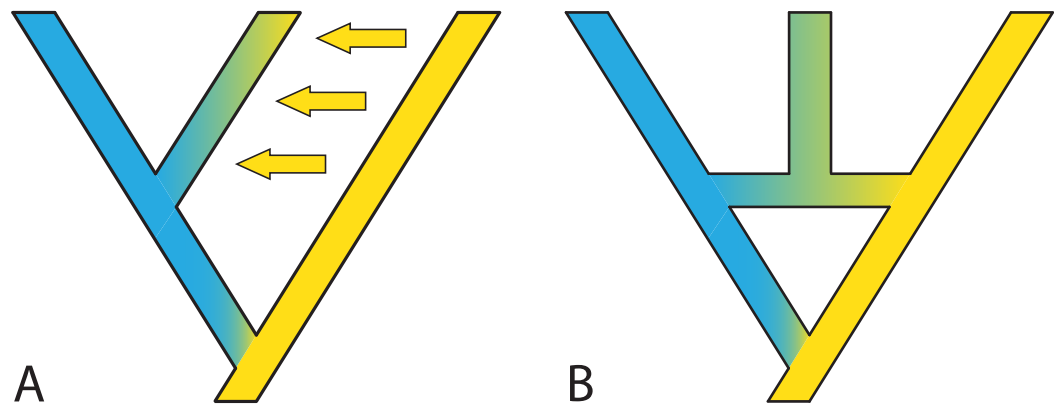


Figure 7 Two possible species phylogenies producing similar observations at present time. (A) species tree with gene flow. (B) Species network with homoploid hybridization. Distinguishing two such scenarios usually requires simulations and comparison of observed and expected summary statistics.

Full-size [DOI: 10.7717/peerj.6399/fig-7](https://doi.org/10.7717/peerj.6399/fig-7)

data as containing genetic information from extinct species, but they do not model the full species phylogeny: such a generalized approach seems a promising avenue to explore.

Polyploidy and the challenges of analyzing gene duplication and loss

The MSC model describes well allelic lineages and the mutations they accumulate (Fig. 3; Degnan & Rosenberg, 2009; Liu, Xi & Davis, 2015). The simple MSC model is challenging to apply to evolutionary events in which the evolving entities (genes or paralogs) duplicate and occasionally go extinct during the evolutionary history of the populations/species and thus cannot be sampled in contemporary population or species. Estimating the existence and number of these “ghost” lineages remains challenging. For example, how can we detect duplication events if one of the duplicated loci is lost in descendant lineages? In the case of polyploidy, two (or more) genomes having separate evolutionary histories end up together in a single individual. What consequences for evolutionary history do genomic conflicts and dosage variation in gene expression impose? Polyploidy also raises technical issues, such as whether or not homoeologous sequences are recovered in standard genomic surveys.

The complication that GDL brings to the inference of species phylogenies has long been recognized (Fitch, 1970). It is therefore surprising that practical solutions to the problem of GDL are almost non-existent, with empirical examples usually based on ad hoc methods and deductions. Ancient duplications where most additional copies are retained in descendent species can be fairly easy to diagnose based on phylogeny (Oxelman et al., 2004; Pfeil et al., 2004). However, resolving duplications becomes more difficult when copy number changes quickly (Ashfield et al., 2012), or when duplications are recent and copy loss is complete or nearly so, thus returning the locus to a single-copy state (Ramadugu et al., 2013). In the latter case, the phylogenetic pattern can mimic that of ILS and become indistinguishable from it (Sousa et al., 2017), generally leaving no trace at all of the loss.

Why is GDL so challenging to implement in theory? The topological and coalescent-time similarities between ILS and GDL complicates extending the MSC to

include both processes, unless copy number exceeds one in at least some samples (Fig. 4). Assuming that allelic and homoeologous variation is not confused with the copy number of independently duplicated genes, at the very least, duplicated genes could be handled as independent loci with missing data for some samples with MSC inference. When copy loss is complete, or when the duplication is so recent so as to conflate allelic versus copy variation, these GDL loci have little effect on species phylogeny inference and divergence times, especially if the algorithms used employ averages over coalescence times or other parameters across many gene trees (Liu et al., 2009; Sousa et al., 2017). At high proportions, though, they may cause serious issues for phylogenetic reconstruction, because the unexpected positions of gene duplications in a species phylogeny, coupled with random copy loss, means that no specific pattern is expected among the affected gene trees (Fig. 4). This scenario contrasts with the retention of ancestral polymorphisms, where we know that branches in short species phylogenies (in coalescent units) are the cause (Rosenberg & Nordborg, 2002). Thus, we expect deeply coalescing lineages to occur in specific parts of a species phylogeny with a limited number of topological outcomes and branch lengths limited by effective population size, which is not the case for duplicated genes. A recent approach to identifying genes that are single copy, but have nonetheless been affected by GDL, was made using the genomic location of the loci (Sousa et al., 2017), and could prove useful for distinguishing GDL and ILS.

Recombination

All existing methods for coalescent estimation of species trees and networks make two important assumptions, namely that (1) there is free recombination between loci, and (2) there is no recombination within a locus. These two assumptions address a key concept distinguishing MSC models from concatenation or supermatrix models: it is the conditional independence of loci, mediated by recombination between loci, and not the ability to address ILS or discordance among genes per se. Moving forward, three important questions to address are: (1) How robust are methods to the presence of recombination within loci and/or to the violation of independence among loci? (2) How should we model recombination within the species phylogeny inference framework? and (3) How do we detect it and differentiate recombination-free loci?

Researchers have started to examine the first question and found a detectable effect of recombination only under extreme levels of ILS and gene tree heterogeneity (e.g., Lanier & Knowles, 2012). However, more analyses and studies are still needed to explore a wider range of factors and parameters that could affect phylogenetic inference when the assumption of recombination-free loci is violated (e.g., Li et al., 2018). For answering the second question, one approach involves combining the multispecies coalescent with hidden Markov models (e.g., Hobolth et al., 2007). These methods suffer from the “state explosion problem”, where individual states are needed for the different coalescent histories, and they increase rapidly with the number of taxa in the dataset, making them infeasible except for very small (~4 taxa) datasets. New methods that scale to larger datasets are needed if such approaches are to be useful in practice. A different direction

is to devise novel methods for inferring species phylogeny while assuming that the genealogies of the individual loci could take the form of an ancestral recombination graph (ARG; [Griffiths & Marjoram, 1996](#); [Siepel, 2009](#); [Rasmussen et al., 2014](#)).

Extending these approaches to address recombination would require the development of new models that significantly extend the multispecies coalescent to account for ARGs within the branches of a species phylogeny. For two-step species tree methods, this entails developing new methods that infer ARGs for the individual loci and methods that infer species phylogenies from collections of ARGs. For single-step Bayesian methods, novel developments are needed to sample species phylogenies, locus-specific ARGs, and their related parameters. It will also be important to better understand the conditions under which ignoring recombination will still yield reasonable estimates of phylogeny. Extending the theory to accommodate ARGs may be of intrinsic interest, but if the parameter space in which recombination is relevant is very small, then practitioners may be able to ignore recombination.

Species concepts and delimitation

Coalescent methods have played an important role in the development and critical evaluation of species delimitation methods because they provide hypotheses for species boundaries based on genetic data and be integrated with phenotypic data (e.g., [Solís-Lemus, Knowles & Ané, 2015](#)). Irrespective of traditional species concepts, it is essential that the entities at the tips of the species tree do not violate the assumptions of the MSC, wherein species are defined mathematically (e.g., [Rannala & Yang, 2003](#), [Degnan & Rosenberg, 2009](#)): the branches of the species tree constitute species or populations that do not exchange genes. However, the MSC model also carries strict assumptions about the divergence process if the delimited units are to be interpreted as species. Specifically, it is important to emphasize that in the “standard” MSC model, these species represent populations that, immediately after divergence, no longer experience gene flow. Therefore, the species of the MSC model do not necessarily correspond with species as a taxonomic rank, defined by traditional species concepts ([Heled & Drummond, 2010](#)): “MSC” species could simply be populations by other criteria, so long as they have ceased to exchange genes, even for a short period of time. In other words, a species tree built under the MSC might then be interpreted as a depiction of the history of the barriers to gene flow among diverging structured populations and this may be particularly true when there is dense spatial and genomic sampling across individuals ([Sukumaran & Knowles, 2017](#)). Therefore, in those species tree methods requiring a priori assignments of individuals to species, such assignments may strongly influence the inferred species phylogeny, in the same way that hybridization will have serious consequences on an estimated species phylogeny ([Leaché et al., 2014](#)).

Recently, several MSC-based methods that have the ability to simultaneously perform species delimitation and estimate the species phylogenies have been developed and implemented (e.g., [Yang & Rannala, 2014](#); [Jones, Aydin & Oxelman, 2015](#); [Jones, 2017](#)). These methods seem to consistently recover the correct number of “MSC species” given the assumptions of the model. However, it is probable that the assumption of no gene

flow between the descendant populations is often violated and that most reproductive isolation processes are gradual or episodic rather than sudden and permanent (e.g., [Rosindell et al., 2010](#)). There is thus need for methods that perform simultaneous species phylogeny estimation and assignment of individuals to species while considering the limitations of the MSC ([Jones, Aydin & Oxelman, 2015](#)).

If one prefers a species concept that affirms that most recently diverged populations are necessarily reproductively isolated, current methods will overestimate the number of species as defined by traditional species concepts and will likely reveal instead intraspecific population structure ([Sukumaran & Knowles, 2017](#)). [Toprak et al. \(2016\)](#) used DISSECT ([Jones, Aydin & Oxelman, 2015](#)) but also employed checks as to the integrity of various hypotheses of species boundaries suggested by the data. From a computational point of view, any species delimitation method will need an operational definition of species. Therefore, a possible development of MSC-based species delimitation methods could be allowing migration and assuming that speciation is complete when a certain proportion of the migrations is reached or when the migration rate is sufficiently low. However, this solution will not be suited for the protracted speciation model because other kinds of information besides the movement of genes will still be needed to identify when a clade becomes reproductively isolated. Possibly the best way to avoid confusion is to restrict the word “species” to taxonomy and base it on multiple sources of information which are synthesized in an integrative fashion ([Dayrat, 2005](#); [Will, Mishler & Wheeler, 2005](#); [Bacon et al., 2012](#); [Solís-Lemus, Knowles & Ané, 2015](#)), and refer to the reproductively isolated units of MSC analysis as “MSC units” or “MSC taxa”.

MODELS AT THE INTERSECTION OF PHYLOGENOMICS AND MACROEVOLUTION

At the intersection of phylogenomics and macroevolution, phylogeneticists aim at shedding light on how patterns of organismal diversity have been generated and maintained through time. The former focuses primarily on building phylogenies, whereas the latter uses them to study the tempo and mode of diversification over time. In many important respects, these two sub-disciplines have remained distinct and non-communicative. On the one hand, phylogenomics and phylogeography have not exhaustively aimed to address the type of questions—related to diversification and trait evolution—that macroevolution focuses on. On the other hand, macroevolution ignores many kinds of complexities inherent to the phylogeny building process that phylogenomics has recently begun to address.

Macroevolutionary models focus on long-term processes, in terms of both species richness and phenotypic diversity. They rely on two types of models: birth-death models of diversification aimed at understanding how and why speciation and extinction rates vary through time and across lineages ([Hey, 1992](#); [Nee, Mooers & Harvey, 1992](#); see [Stadler, 2013](#) and [Morlon, 2014](#) for review) and models of trait evolution aimed at understanding the mode and tempo of phenotypic evolution ([Felsenstein, 1973](#); see [Pennell & Harmon, 2013](#) and [Manceau, Lambert & Morlon, 2017](#) for reviews). These models are typically constructed at the level of species, ignoring the populations or individuals that

constitute these species (but see [Manceau, Lambert & Morlon, 2015](#) and [Rosindell, Harmon & Etienne, 2015](#) for exceptions). As a consequence, microevolutionary processes, such as coalescence, have informed phylogenetic methods for building phylogenies more so than have macroevolutionary methods that use them. For example, the most widely used phylogenetic dating methods generally do not acknowledge the critical distinction between speciation times, which are usually of primary interest, and coalescence times, which are often assumed to represent speciation times but in fact represent events older than the divergence of the species concerned ([Edwards & Beerli, 2000](#); [dos Reis, Donoghue & Yang, 2016](#); [Angelis & dos Reis, 2015](#)). In addition, macroevolutionary models are fit to species phylogenies (diversification models) or a combination of species phylogenies and phenotypic data (trait evolution models), most often assuming that evolution is best represented by a species tree, not a network (but see [Jhwueng & O'Meara \(2015\)](#); [Bastide et al. \(2018\)](#); [Solís-Lemus, Bastide & Ané \(2017\)](#) for models of trait evolution on networks), and that the species phylogeny is known. Nearly all models that use phylogenies to study character evolution assume a single underlying species phylogeny on which characters evolve. But it has become evident recently that different characters in principle have different phylogenies, for the same reason that genes themselves might have different phylogenies ([Hahn & Nakhleh, 2016](#)). Analyzing incongruences between character evolution inferred from the species tree versus from gene trees that are more directly linked to the character under study would provide a refined understanding of character evolution. Recent work on the phylogeny of quantitative characters may be helpful in this endeavor ([Felsenstein, 2012](#)).

Developing research projects that integrate the heterogeneity inherent in phylogenomics and macroevolution will bring important new insights into the evolutionary process. For example, developing diversification and phenotypic evolution models to be fit to networks rather than dichotomous trees will allow estimates of rates of hybrid speciation and phenotypic evolution as well as a better understanding of factors influencing such rates ([Bastide et al., 2018](#)). Embracing genetic heterogeneity and the incongruence between gene trees and species phylogenies when applying macroevolutionary models could help us to better understand how speciation proceeds, and also to analyze the coupling between genetic and phenotypic evolution (e.g., is phenotypic convergence coupled or not with genetic convergence in relevant genes?). Developing macroevolutionary models accounting for within-species heterogeneity linked to biogeography could help us understand how biogeographic structuring influences speciation, extinction, and phenotypic evolution.

More generally, evolutionary biologists have yet to thoroughly explore the type of new questions that we are going to be able to address if we are given genomic data at the tips of all species from a phylogeny. Such data could allow us to gain an integrative understanding of three fundamental aspects of evolution: evolution at the molecular level, at the phenotypic level, and at the clade level, as well as the links among them. Are rates of evolution at these three levels correlated? If so, how? Do features of genomes or of genome evolution, such as quantity of transposable elements, substitution rates, number of

gene duplications, influence rates of diversification and phenotypic evolution? Clearly, we are only at the beginning of exploring these new possibilities.

Mapping trait evolution on heterogeneous genomic datasets

Mapping the genomic basis of phenotypic traits is a major trend in evolutionary biology today ([Elmer & Meyer, 2011](#); [Hoban et al., 2016](#)). Such mapping can be conducted in the context of populations of a single species or, increasingly, via comparisons of species on a phylogeny (e.g., [Hiller et al., 2012](#); [Marcovitz, Jia & Bejerano, 2016](#)). Phylogenetic genome-wide association (“PhyloGWAS”) methods identify genomic features in coding or non-coding DNA that exhibit unusual patterns of evolution on branches concerned with repeated evolution of phenotypes, thereby drawing connections between the genomic and phenotypic levels ([Pease et al., 2016](#)). Such phylogenomic mapping usually assumes a single phylogeny, the species phylogeny, as a framework for analysis, and therefore ignores genomic heterogeneity. To make PhyloGWAS mapping most efficient it might be more appropriate to use the local topology in the genome for inference and estimation of ancestral states. Estimating genotype-phenotype associations solely on the species phylogeny might yield misleading results regarding the origin and evolution of phenotypic traits ([Hahn & Nakhleh, 2016](#); [Mendes, Hahn & Hahn, 2016](#); [Guerrero & Hahn, 2018](#)). Heterogeneity across gene histories has been traditionally considered as “biological noise” when using comparative genomics to map traits, but of course such heterogeneity is the focus of gene mapping efforts at lower taxonomic levels. A recently proposed application of the MSC for quantitative traits, accounting for genealogical heterogeneity improves downstream estimates of mean trait values, phylogenetic signal, and evolutionary rates of traits ([Mendes et al., 2018](#)). Furthermore, genome-wide or gene-specific selective sweeps associated with the evolution of a particular phenotypic trait are a major source of genetic heterogeneity among closely related populations or species, and can be captured using outlier statistics, such as F_{st} or D_{xy} ([Pease et al., 2016](#)). Such selective sweep mapping of genes with large phenotypic effect can now be accomplished with high resolution and precision in genomically poorly studied organisms ([Lamichhaney et al., 2015](#)). Apart from providing valuable knowledge on the genetic basis of trait diversification, such data are providing increasing support to the fact that cases of genetic heterogeneity can be profitably used in the effort to understand and resolve evolutionary history, rather than considering it “biological noise.” Such thinking needs to be incorporated into comparative genomics more frequently (e.g., [Mendes et al., 2018](#)).

Tree-free methods of character evolution

We have seen that incorporating phylogenetic heterogeneity is a challenge for macroevolutionary models of character evolution. At the other end of the spectrum are a class of methods (so called “tree-free methods”) that attempt to draw inferences and principles about trait evolution without assuming a particular phylogeny. The common situation when analyzing character or trait data correlated by a phylogeny is to assume a stochastic process for the trait, commonly a variation of the Brownian motion

(BM; [Felsenstein, 1985](#)) or Ornstein-Uhlenbeck (OU; [Hansen, 1997](#)) processes. Then, using the estimated phylogeny and measured trait data for each species, the parameters of various evolutionary processes—trait variation, patterns and rates of change, etc.—are estimated, often using maximum-likelihood or Bayesian approaches (see [Pennell & Harmon, 2013](#) and [Manceau, Lambert & Morlon, 2017](#) for reviews). However, given the various logistical and technical challenges of inferring robust phylogenies, exploring tree-free methods might represent a useful mechanism for guiding the study of character evolution for certain groups.

Tree-free comparative methods work by integrating over the space of trees (under a given branching process model). For example, under a pure birth model and with enough tip measurements, the optimum value of the OU process can be estimated as the sample average ([Bartoszek & Sagitov, 2015a](#)). Similar results have now been derived for other models of tree growth that include extinction ([Adamczak & Miloś, 2014; 2015; Ané, Ho & Roch, 2017](#)). Similarly, the rate of adaptation under the OU process, often modeled as the stationary variance—the ratio of the squared “rate of evolution” (sigma parameter in the OU model) and twice the “rate of adaptation” (the alpha parameter) can be estimated as the sample variance ([Bartoszek & Sagitov, 2015a](#)). Teasing sigma and alpha apart, however, requires a tree. The key parameter of the BM model, the rate of evolution, is similarly estimable directly from the trait sample ([Bartoszek & Sagitov, 2015b; Crawford & Suchard, 2013](#)), whereas the root state cannot be consistently estimated without a tree ([Ané, 2008; Sagitov & Bartoszek, 2012](#)). In addition to providing tree-free estimators of some model parameters, the studies mentioned above also derived Central Limit Theorems that allow computing confidence intervals around these point estimates as well as the sample sizes needed to obtain reliable estimates.

Extinct and unsampled species

A notable case when phylogenomics and macroevolution do meet is in the treatment of extinct or unsampled species in phylogenetic reconstruction and dating. Despite the avalanche of genomic data for an increasing number of species, we still lack sequence data for most species, making it difficult to place them in a phylogeny. Some researchers (e.g., [Jetz et al., 2012; Tonini et al., 2016](#)) have used a combination of consensus trees and current taxonomic classifications to impute the phylogenetic relationships of unsampled species. In this case, polytomies are often resolved by using distributions of branching times obtained from macroevolutionary birth-death models ([Kuhn, Mooers & Thomas, 2011](#)). While we appreciate the value of these approaches given the real logistic difficulties researchers face as they attempt to obtain samples from around the globe and that methods are now easily applicable, such approaches have generated well-founded concerns about biases in our inferences ([Davies et al., 2012; Rabosky, 2015](#)) and the extent of these biases remain largely unknown ([Pennell, FitzJohn & Cornwell, 2016](#)). Nonetheless, there are methods that can objectively alleviate issues arising from using incomplete phylogenies in downstream phylogenetic comparative analyses. For example, recent results using conditioned birth-death processes (e.g., [Gernhard, 2008a, 2008b; Sagitov & Bartoszek, 2012](#)) show that under constant rate processes the size

of the clade contributes information on the height of the tree and also on the coalescence times. Such results can be used to improve the calibration and node dating of the phylogeny when some species are not sampled. One would expect that ignoring the non-sequenced species would incur a bias resulting in shorter tree heights, because less time is usually required to generate fewer tips. Conditioned branching process models can help alleviate this bias. Also, macroevolutionary birth-death models are used as branching process priors in Bayesian molecular dating. The availability of likelihood expressions for incompletely sampled phylogenies ([Stadler, 2009](#); [Stadler & Steel, 2012](#); [Morlon, Parsons & Plotkin, 2011](#)) thus allow to date phylogenies while accounting for the fact that we have observed only a certain fraction of unsampled species.

An important source of genetic data for extinct and unsampled species comes from ancient DNA. Current methods now allow sequencing ancestral DNA and incorporating it into phylogenomic, phylogeographic, and population genetic analyses (reviewed by [Leonardi et al., 2017](#)), and even obtaining wholegenomes from various samples dating back up to a few thousand years (e.g., [Lynch et al., 2015](#); [Cloutier et al., 2018b](#)). The paleontological and paleobiological records represent another source of extinct species for which genomic data is out of reach. Having the opportunity of integrating the temporal and phylogenetic information of extinct data provided by fossils into phylogenomic estimates of extant diversity is a crucial task toward a complete ToL (reviewed by [Hunt & Slater, 2016](#)). Besides the possibility of placing fossil data into phylogenetic trees, a growing field of development is that of time-scaling phylogenies using fossil data. Traditionally, node-dating methods have used fossils to bound the minimum age of nodes at the base of extant clades in which fossils likely belong based on morphological comparisons (reviewed by [Ksepka et al., 2011](#)) or estimated fossil ages have been used to post-hoc-scale branch lengths ([Sanderson, 2002](#); [Smith & O'Meara, 2012](#); [Bapst, 2014](#)). However, there is an increasing interest in simultaneously performing phylogenetic inference and time-scaling phylogenies ([Pyron, 2011](#); [Ronquist et al., 2012](#); [Heath, Huelsenbeck & Stadler, 2014](#)). These methods, known as tip-dating, have the advantage of accommodating phylogenetic uncertainty implemented in a Bayesian model-based framework ([Hunt & Slater, 2016](#)) and have been implemented in several groups spanning a wide temporal scale (e.g., [Pyron, 2011](#); [Pyron & Burbrink, 2012](#); [Ronquist et al., 2012](#); [Heath, Huelsenbeck & Stadler, 2014](#); [Slater, 2015](#); [Larson-Johnson, 2016](#); [Matzke & Wright, 2016](#)). Lastly, these phylogenies containing both extinct and extant diversity serve as framework to conduct robust diversification and trait evolution studies ([Hunt, 2013](#); [Pennell & Harmon, 2013](#); [Pyron, 2015](#)).

BUILDING, UPDATING AND SUSTAINING THE TREE OF LIFE

Scalability challenges

Inferring the phylogeny of all living organisms represents a different challenge than inferring the relationships of just a few terminals; often the scale at which new methods are developed and tested is on this latter scale. For instance, for eukaryotes alone, recent conservative estimates indicate that there are ~8.7 million species on Earth and only 9–14% of them have been formally described ([Mora et al., 2011](#)). Furthermore,

out of 2.6 million taxa currently represented in the Open Tree of Life (<https://tree.opentreeoflife.org>; *Hinchliff et al., 2015*), only ~55,000 were gathered from hard-data phylogenies, whereas phylogenetic affinities of the rest were inferred from current taxonomic classifications (*McTavish et al., 2015, 2017*). These observations suggest that the vast majority of taxa on Earth still await formal taxonomic description and placement in the ToL (*Mora et al., 2011; McTavish et al., 2017*). As mentioned above, one common challenge that phylogeneticists encounter is the difficulty in accessing samples from rare, endangered, or extinct taxa, particularly in countries where collecting and exporting is difficult. Recent genomic techniques now allow successful results in obtaining valuable ancient DNA data from museum specimens (e.g., *Staats et al., 2013; Hykin, Bi & McGuire, 2015; McCormack, Tsai & Faircloth, 2016; McCormack et al., 2017; Ruane & Austin, 2017*), and here, we advocate for routine use of these resources to improve taxon sampling, enhance research in phylogenomics and phylogeography, and increase awareness and usefulness of natural history collections.

Despite the great increase in the generation of genomic data across organisms, we are often limited to fast and efficient but simpler, less realistic phylogenetic methods and assumptions, such as IQ-Tree (*Nguyen et al., 2015*), PhyloBayes (*Lartillot et al., 2013*), and ExaML (*Kozlov, Aberer & Stamatakis, 2015*), to deal with large, heterogeneous datasets. The main downside of these options is that they are not a full coalescent framework and thus far rely on data concatenation. Fully coalescent methods, such as the popular phylogenetic software program *BEAST (*Heled & Drummond, 2010*), are not capable of dealing with more than a few hundred taxa and some dozen loci at a time for a common analysis, and only recently the release of StarBeast2 allows for the use of thousands of loci for tens of taxa (*Ogilvie, Bouckaert & Drummond, 2017*). To tackle this problem, we encourage the continuing development of methods that are fully scalable and ideally only increase analytical time linearly rather than exponentially with the number of taxa and loci. Phylogenetic methods should also be fully parallelizable (in order to run natively in computer clusters) and contain checkpoints, i.e., be able to resume the analyses from the latest logged file in case an analysis crashes or the user wishes to evaluate partial results. Another point of possible improvement is in dealing with new sequences to be added to a previously large dataset: should the analysis start from scratch, or could there be substantial time gains by letting those sequences find their placement in the phylogeny ‘on the fly’ (e.g., *Siu-Ting et al., 2014*).

Large scale phylogenies should ideally be based on the best (or most comprehensive) available datasets in terms of taxonomic and molecular sampling and be constructed from the data itself. However, even supermatrix inference conducted under a single analysis can add bias on tree heights and coalescence times when performed across unbalanced sampled clades (a very common case for species-rich clades or understudied taxa), and therefore affect downstream analyses that rely on these parameters (e.g., biogeography, trait evolution, diversification rates). Computing optimally populated datasets that combine the largest number of taxa and loci simultaneously is a complex mathematical problem, but recent approaches (e.g., SUPERSMART—*Antonelli et al., 2017*; PyPHLAWD—*Smith & Brown, 2018*) attempt to overcome it objectively,

such as applying the knapsack problem to phylogenetics by packing the optimal choice of species and suitable alignments into a minimally sparse supermatrix.

Community initiatives

Building the ToL is a grand challenge in molecular phylogenetics, and one that cannot be accomplished by a single person or institution's efforts. Several initiatives have been developed in recent years to coordinate efforts and provide the research community with synthetic information. A prominent project is the Open Tree of Life (<https://tree.opentreeoflife.org/>; [Hinchliff et al., 2015](#)). This project provides a synthesis of previously published phylogenies merged through supertree and other grafting methods. One issue faced by the initiative is that it relies on authors uploading their phylogenetic trees to open data repositories, such as Dryad Data Repository (<http://datadryad.org/pages/organization>; [Vision, 2010](#)) or TreeBase ([Sanderson et al., 1994](#); [Piel et al., 2009](#)), which at least until recently only occurred in about 17% of cases ([Drew, 2013](#)). Substantial curatorial efforts are also critical to facilitate reusability of deposited trees ([McTavish et al., 2015](#)). A different approach was taken by [Antonelli et al. \(2017\)](#), who developed a framework for continuously inferring time-calibrated large phylogenies from raw sequence data deposited in GenBank ([Clark et al., 2016](#)) in a multi-step method. Similarly, various tools have been developed to make information contained in the ToL available for the general public (e.g., [Rosindell & Harmon, 2012](#); [Harmon et al., 2013](#)).

Mapping the Tree of Life

While progress has been made in mapping species distributions at the large scale aiming for improved conservation practices (e.g., the Map of Life collaborative project; <https://mol.org/>), most initiatives do not map the tips of phylogenetic trees directly onto geographic space, and therefore are limited by current taxonomic knowledge. As spatial variation in biodiversity results from interactions between evolutionary history and environmental factors, explicit connections between the tips of the ToL and geographic ranges will greatly improve biogeographic inferences ([Quintero et al., 2015](#)) and our understanding of biodiversity patterns and future trends. Advances in mapping the ToL through earth history using genomic-based phylogenetic inferences over broad scales and explicit spatial models (e.g., geophylogenies and continuous diffusion models; [Kidd, 2010](#)) depend directly on locality data that should be made available in raw and ready-to-use formats. Data sharing policies for associated data, such as geographic coordinates and voucher information, is not well established among journals. We argue that editorial boards should try as best as possible to establish data policies that value and encourage the deposit of geographic data associated to vouchered specimens and other associated information available for future reference.

Best practices for building the Tree of Life

Data must be well curated in databases and publicly available

As we are now in the era of big data in biological sciences, adequate reproducibility must be a fundamental endeavor of biodiversity research. Therefore, data publication in open

access repositories represents a powerful tool that not only ensures long-term storage and public availability for future research, but also serves as a vehicle for clarifying intellectual rights and scientific merits (Costello & Wieczorek, 2014). Biocuration, the activity of organizing, representing, and making biological information accessible to both biologists and bioinformaticians, has now become an important consideration in building, updating, and sustaining the ToL (McTavish et al., 2017). The exponential growth in the amount of genomic scale data and the increased dependence on the availability of each other's' data to answer complex biological questions means that there is a need for improved data management, analysis, and accessibility. GenBank has been the main open access repository for annotated collections of publicly available molecular data. Although the data stored in this database usually lists information such as organism of origin and publication details, the utility of molecular data in this database to answer multiple biological questions, such as biogeographic patterns of biodiversity, is often hampered by lack of associated information such as collection locality (Scotch et al., 2011; Gratton et al., 2017) or attachment to a specific voucher specimen. Moreover, recent surveys have shown that fewer than 20% of phylogenetic studies provide access to phylogenetic data (i.e., alignments and phylogenies) and when they do, critical biological information such as complete taxon names is missing (Drew et al., 2013; Magee, May & Moore, 2014).

We propose two urgent actions to advance this key field. First, authors should be encouraged to submit molecular data that is linked to voucher specimens deposited in recognized scientific collection and museums. Second, authors, journals, and curators should encourage all molecular data submitted to include information such as collection locality and details of voucher specimens. In this regard, other global initiatives such as the International Barcode of Life Project (iBOL; <http://www.ibolproject.org>) have had great success linking molecular data with morphological and distributional data. When all the data produced or published are curated to high standards and made accessible as soon as available, biological research will be able to process massive amounts of complex data much more quickly.

Submitting sequence and tree data during publication is now routine. However, making available all analytical methods such as software and code used to process and analyze data is less widely employed by the phylogenetic community. Facilities such as TreeBase, Dryad Digital Repository, and Github (<https://github.com/>) provide a platform for the curated storage of the data and bioinformatic pipelines underlying the scientific literature (see McTavish et al., 2015; 2017). Authors and journals should require all published research to include links to raw data, processed data, and all analytical methods used to produce the results presented. In general, we advocate for following best practices of data management and publication to ensure the quality and utility of phylogenomic data and their associated biological information (Stoltzfus et al., 2012; Drew et al., 2013; Costello & Wieczorek, 2014). In putting together Fig. 2, for example, we found that basic information on a given phylogenomic study, such as the number of species or sequences analyzed, or the total number of base pairs in an alignment, were often not reported or difficult to recover; including such information in easy-to-access

tables prior to article acceptance would greatly facilitate meta-analyses and syntheses as the number of studies grows (Table S1).

The need for adequate curation of analytical tools

In the same way that data must be adequately stored and curated, analytical tools must be available for future use and should guarantee proper reproducibility (Wilson *et al.*, 2014; Darriba, Flouri & Stamatakis, 2018; DeBiasse & Ryan, 2018). One of the reasons behind the dramatic increase in the number of phylogeographic and phylogenetic studies during the last 20 years is the proliferation of software and bioinformatic tools to process and analyze these data. Thanks to these new methods, it is now possible to implement a wide array of theoretical models that sustain the fields of phylogenomics and phylogeography. As stated above, genome-wide data have notoriously increased the necessity to expand our analytical models, ultimately leading to a stronger demand for computing resources (Darriba, Flouri & Stamatakis, 2018). Given their key role in phylogenomic research, it is advisable that software development, documentation, and availability follow the best possible practices (e.g., Leprevost *et al.*, 2014; Wilson *et al.*, 2014; Guang *et al.*, 2016; Darriba, Flouri & Stamatakis, 2018; DeBiasse & Ryan, 2018). Having both data and analytical tools adequately stored and accessible to the public not only will ensure high reproducibility of previous studies, but, more importantly, will facilitate continuing the construction of the ToL (McTavish *et al.*, 2017).

All contributions toward building the Tree of Life must be properly recognized

Some current publishing practices in the scientific community may unintentionally represent hurdles toward the ultimate end of collecting and disseminating phylogenetic data on which to build a ToL. For instance, the increasing need in many countries and communities for publishing high-impact papers understandably often discourages researchers from releasing their data until their studies are complete and have passed the peer-review process. This is partially explained by the heavy emphasis of top journals on unusually novel and flashy findings as compared to those studies that represent more modest, but just as critical, advances in the understanding of the phylogenetic relationships of the groups. Similarly, this urge to publish high-impact papers often impedes adequate long-term studies that could potentially generate a wider variety of basic data. With the cultural emphasis on impact and numbers and rates of publication, in practice there is often a penalty for long-term studies. Our current climate often values novel results produced in the short term. Consequently, as a community, we must reach an equilibrium between short- and long-term scientific production in a way that values both, encouraging high impact studies bringing radical reorganizations of the ToL, without hurting lower impact research and the ongoing search for innovation.

Moreover, because building the ToL is a slow and daunting task, it is important that, as a scientific community, all contributors to the process receive proper recognition for their contributions, thereby keeping motivation high and retaining our best talent. Unfortunately, some contributors, both institutions and roles within them, receive less recognition in this grand task than others. For example, field biologists that obtain basic

Table 2 Challenges in the fields of coalescent-based phylogenomics and implications for unraveling character evolution and the Tree of Life.

| Category | Challenge | Proposed strategy |
|---|---|--|
| Data | Integration and assessment of large amounts of data with heterogeneous phylogenetic signal. | <p>Protocols for marker selection should assess markers' biological relevance and adequacy for the study organism, given the temporal and spatial scales in question, and not only logistical convenience. A posteriori (after data generation) marker selection from whole-genome alignments can be useful to inform these aspects as well as minimize the effects of missing data and varying data quality. Until then, researchers should attempt a higher standardization of markers to facilitate combinatory analyses.</p> <p>To discern true phylogenomic heterogeneity from noise and error as well as to identify violations of the MSC, adequate filtering of large phylogenomic datasets should be conducted based on biological and statistical properties of markers (e.g., analyses of gene-tree outliers and rogue taxa).</p> <p>Further research on filtering methods as well as on their impact on phylogenomic estimation is still required.</p> |
| | Inclusion of additional character types into phylogenomic analyses. | <p>Research efforts focused on the adequate identification and utilization of rare genomic changes other than nucleotide substitutions, such as indels, transpositions, inversions, CNVs, and chromosomal rearrangements. Development of new methods not only to infer phylogenetic hypotheses based on these characters but also to integrate them with more traditional sequence data.</p> |
| Phylogenetic inference models and methods | Analyses of genome-scale data for large numbers of species within a coalescent framework. | <p>Continue the development of models and methods that allow simultaneous gene tree and species tree estimation within a Bayesian framework (e.g., Ogilvie, Bouckaert & Drummond, 2017) for increasingly large and complex datasets.</p> <p>For the time being, two-step methods, particularly those based on biological models and permitting statistical tests of topologies in a likelihood framework (e.g., Liu, Yu & Edwards, 2010), are useful tools to incorporate coalescent information into species tree inference.</p> |
| | Detection and incorporation of violations of the MSC into phylogenomic inferences. | <p>Extensions of the MSC should seek the inference of reticulate evolutionary histories (i.e., multispecies network coalescent; Yu et al., 2014; Wen et al., 2016) by simultaneously incorporating violations of the MSC (e.g., Wen & Nakhleh, 2018 (reticulation and ILS); Jones, Sagitov & Oxelman, 2013 (allopolyploidy)). Inference methods dealing with GDL and recombination are of high priority.</p> <p>Further development of conceptual approaches aimed at detecting and quantifying different underlying biological processes of phylogenetic history (e.g., Jones, 2018 (ILS and migration), Blischak et al., 2018 (hybridization), Hibbins & Hahn, 2018 (direction and timing of introgression), Sousa et al., 2017 (ILS and GDL), Li et al., 2018 (recombination rates)). As proposed by Degnan (2018), using the MSC as a null model within a model selection approach can be a powerful tool to identify violations of the MSC and to deepen our understanding of the biological consequences of these processes.</p> |

Table 2 (continued).

| Category | Challenge | Proposed strategy |
|--|--|---|
| Models that integrate phylogenomics and comparative analyses | Integrating different phylogenetic signals into comparative analyses. | Methods and models should attempt to incorporate gene tree incongruence into macroevolutionary models of character evolution. Similarly, integrative studies aiming at unraveling character evolution at the molecular, phenotypic, and clade levels. |
| | Understanding the genomic bases of character evolution in species trees vs. gene trees. | Methods that estimate phenotype-genotype associations incorporating heterogeneity across gene trees or that at least take into account differential state probabilities stemming from gene tree discordance (e.g., Guerrero & Hahn, 2018). Similarly, extensions of the MSC for quantitative traits that take into account genealogical heterogeneity represent a promising avenue for research and implementation (e.g., Mendes et al., 2018). |
| Best practices for building the ToL | Increasing the number of species represented in the ToL while ensuring reproducibility and encouraging community participation | Natural history museums must be central players for providing and analyzing genome-scale data. Genetic resources and specimen collections are fundamental for allowing the acquisition of data for extinct and poorly accessible species. Open access community initiatives must continue to be relevant repositories of the ToL. Adequate methods for curation of data and analytical tools must continue to be a high priority. |

natural history information and specimens used for building the ToL ([Suarez & Tutsui, 2004](#)), and the natural history museums that house those specimens, are often not recognized sufficiently. As a community, we have been following a trend in which, perhaps inadvertently, we do not value as much the production of basic biological and natural history data. This can certainly be recognized in our national funding practices, which often do not support basic taxonomic or natural history fieldwork at the expense of flashier end-uses of biological specimens. Specimens are the foundation of most phylogenomic and phylogeographic studies, and we should find standard mechanisms not only to acknowledge, but also to encourage the production of these data in an integrative framework. It is time to strengthen those initiatives aimed at recognizing scientific production beyond citations of peer-reviewed literature (e.g., ORCID; <https://orcid.org>) by giving also credit to the production and impact of basic biology datasets and collected specimens. Providing credit for depositing and generating data by tracking, for example, number of access and downloads or number of studies using genetic data associated to specimens, could represent a formal recognition of the importance of producing and sharing basic biological data could help bridge the gap between naturalists, taxonomists, empiricists, and mathematicians invested on the study of life history.

It will be exciting to have objective estimates that allow tracking the direct and indirect impact of how these data and samples are being used. We are confident that such initiatives will highlight the importance of continuing field- and museum-based research in various fields of biological research ([Buerki & Baker, 2016](#)). Furthermore, such cultural shifts will undoubtedly encourage discerning young minds to embrace basic biological research in their academic endeavors, rather than embracing more lucrative and societally appreciated applied fields.

CONCLUSIONS

In this perspective, we have attempted to cover ground in the vast arena of issues facing modern phylogenomics today. In [Table 2](#), we summarize some of the most pressing challenges that the field of phylogenomics is experiencing as we use coalescent-based methods toward building the ToL. We have seen how genome-scale phylogenomics, currently on a strong footing as a result of the MSC, is increasingly improved by models that recognize reticulate processes, such as recombination and introgression. In contrast, macroevolutionary models that use phylogenies have yet to embrace the heterogeneity that currently drives many theoretical innovations in phylogenetic reconstruction itself. We have emphasized the need for the phylogenomics community to embrace high standards of data quality, curation and accessibility in its long-term pursuit of the ToL. Such a grand mission requires value and recognition placed not only on the end products of the process, such as publications and trees, but also on the natural history specimens on which phylogenies are based and which are cared for by the community of natural history museums. Building the ToL will require contributions from all sectors of biological and related sciences—from field biology to theory and everything in between—and robust cyberinfrastructures to integrate these diverse and increasingly massive data streams.

ACKNOWLEDGEMENTS

This paper is a product of the ‘Origin of Biodiversity Workshop’ organized by Chalmers University of Technology and the University of Gothenburg, under the auspices of the Gothenburg Centre for Advanced Studies (GoCAS). We are particularly grateful to the GoCAS organizers and facilitators, in particular Karin Hårding, Mattias Marklund, Bernt Wennberg, Sandra Johansson, and Lotta Fernström. We thank Johnathan Clark, Alison Cloutier, Phil Grayson, Kathrin Näpflin, Flavia Termignoni, Jonathan Schmitt, Simon Sin, João Tonini, and Pengcheng Wang for help compiling [Table S1](#). Thomas Couvreur, Tobias Andermann, Prosanta Chakrabarty, Alex Pyron, Claire Morgan, Chris Creevey, and one anonymous reviewer provided useful comments that improved the contents of this manuscript.

ADDITIONAL INFORMATION AND DECLARATIONS

Funding

The Gothenburg Center for Advanced Studies (GoCas) workshop ‘Origins of Biodiversity’ was funded by Chalmers University of Technology and the University of Gothenburg. The authors are supported by scholarships or research grants from the following agencies: Swedish Research Council (Bengt Oxelman, Alexandre Antonelli); U.S. National Science Foundation; European Research Council under the European Union’s Seventh Framework Programme (FP/2007-2013, ERC Grant Agreement n. 331024 to Alexandre Antonelli); Swedish Foundation for Strategic Research; Wallenberg Academy Fellowship (Alexandre Antonelli); Conselho Nacional de Desenvolvimento Científico e Tecnológico—CNPq (Fernanda P. Werneck); Partnerships for Enhanced Engagement in Research from the U.S.

National Academy of Sciences (Fernanda P. Werneck); U.S. Agency of International Development—PEER NAS/USAID (Fernanda P. Werneck); and L’Oreal-Unesco For Women in Science Program (Fernanda P. Werneck). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Grant Disclosures

The following grant information was disclosed by the authors:

Chalmers University of Technology.

University of Gothenburg.

Swedish Research Council (Bengt Oxelman, Alexandre Antonelli).

U.S. National Science Foundation.

European Research Council under the European Union’s Seventh Framework Programme: FP/2007-2013, ERC Grant Agreement n. 331024.

Swedish Foundation for Strategic Research.

Wallenberg Academy Fellowship.

Conselho Nacional de Desenvolvimento Científico e Tecnológico—CNPq.

Partnerships for Enhanced Engagement in Research from the U.S. National Academy of Sciences.

U.S. Agency of International Development—PEER NAS/USAID.

L’Oreal-Unesco For Women in Science Program.

Competing Interests

Alexander Schliep and Scott V. Edwards are Academic Editors for PeerJ.

Author Contributions

- Gustavo A. Bravo analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the paper, approved the final draft.
- Alexandre Antonelli analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Christine D. Bacon analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Krzysztof Bartoszek analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Mozes P. K. Blom analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Stella Huynh analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Graham Jones analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the paper, approved the final draft.
- L. Lacey Knowles analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Sangeet Lamichhaney analyzed the data, authored or reviewed drafts of the paper, approved the final draft.

- Thomas Marcussen analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Hélène Morlon analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Luay K. Nakhleh analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Bengt Oxelman conceived and designed the experiments, analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Bernard Pfeil analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Alexander Schliep analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Niklas Wahlberg analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Fernanda P. Werneck analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- John Wiedenhoeft analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Sandi Willows-Munro analyzed the data, authored or reviewed drafts of the paper, approved the final draft.
- Scott V. Edwards conceived and designed the experiments, analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the paper, approved the final draft.

Data Availability

The following information was supplied regarding data availability:

The raw data are available in [Table S1](#).

Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj.6399#supplemental-information>.

REFERENCES

- Aberer AJ, Krompass D, Stamatakis A. 2013.** Pruning rogue taxa improves phylogenetic accuracy: an efficient algorithm and webservice. *Systematic Biology* **62**(1):162–166 DOI [10.1093/sysbio/sys078](https://doi.org/10.1093/sysbio/sys078).
- Adamczak R, Miloś P. 2014.** U-statistics of Ornstein-Uhlenbeck branching particle system. *Journal of Theoretical Probability* **27**(4):1071–1111 DOI [10.1007/s10959-013-0503-2](https://doi.org/10.1007/s10959-013-0503-2).
- Adamczak R, Miloś P. 2015.** CLT for Ornstein-Uhlenbeck branching particle system. *Electronic Journal of Probability* **20**:1–35 DOI [10.1214/EJP.v20-4233](https://doi.org/10.1214/EJP.v20-4233).
- Andermann T, Fernandes AM, Olsson U, Topel M, Pfeil B, Oxelman B, Aleixo A, Faircloth BC, Antonelli A. 2018.** Allele phasing greatly improves the phylogenetic utility of ultraconserved elements. *Systematic Biology* **68**(1):32–46 DOI [10.1093/sysbio/syy039](https://doi.org/10.1093/sysbio/syy039).
- Ané C. 2008.** Analysis of comparative data with hierarchical autocorrelation. *Annals of Applied Statistics* **2**(3):1078–1102 DOI [10.1214/08-AOAS173](https://doi.org/10.1214/08-AOAS173).

- Ané C, Ho LST, Roch S. 2017. Phase transition on the convergence rate of parameter estimation under an Ornstein-Uhlenbeck diffusion on a tree. *Journal of Mathematical Biology* 74(1–2):355–385 DOI 10.1007/s00285-016-1029-x.
- Angelis K, Dos Reis M. 2015. The impact of ancestral population size and incomplete lineage sorting on Bayesian estimation of species divergence times. *Current Zoology* 61(5):874–885 DOI 10.1093/czoolo/61.5.874.
- Antonelli A, Hettling H, Condamine FL, Vos K, Nielsson RH, Sanderson J, Sauquet H, Scharn R, Silvestro D, Töpel M, Bacon CD, Oxelman B, Vos RA. 2017. Towards a self-updating platform for estimating rates of speciation and migration, ages, and relationships of taxa. *Systematic Biology* 66:152–166 DOI 10.1093/sysbio/syw066.
- Ashfield T, Egan AN, Pfeil BE, Chen NWG, Podicheti R, Ratnaparkhe MB, Ameline-Torregrosa C, Denny R, Cannon S, Doyle JJ, Geffroy V, Roe BA, Saghai-Marroof MA, Young ND, Innes RW. 2012. Evolution of a complex disease resistance gene cluster in diploid *Phaseolus* and tetraploid *Glycine*. *Plant Physiology* 159(1):336–354 DOI 10.1104/pp.112.195040.
- Ashkenazy H, Cohen O, Pupko T, Huchon D. 2014. Indel reliability in indel-based phylogenetic inference. *Genome Biology and Evolution* 6(12):3199–3209 DOI 10.1093/gbe/evu252.
- Avise JC, Arnold J, Ball RM, Bermingham E, Lamb T, Neigel JE, Reeb CA, Saunders NC. 1987. Intraspecific phylogeography: The mitochondrial DNA bridge between population genetics and systematics. *Annual Review of Ecology and Systematics* 18:489–522 DOI 10.1146/annurev.es.18.110187.002421.
- Bacon CD, McKenna MJ, Simmons MP, Wagner WL. 2012. Evaluating multiple criteria for species delimitation: an empirical example using Hawaiian palms (Arecaceae: *Pritchardia*). *BMC Evolutionary Biology* 12(1):12–23 DOI 10.1186/1471-2148-12-23.
- Baker AJ, Haddrath O, McPherson JD, Cloutier A. 2014. Genomic support for a Moa-Tinamou clade and adaptive morphological convergence in flightless ratites. *Molecular Biology and Evolution* 31(7):1686–1696 DOI 10.1093/molbev/msu153.
- Bapst DW. 2014. Assessing the effect of time-scaling methods on phylogeny-based analyses in the fossil record. *Paleobiology* 40(03):331–351 DOI 10.1666/13033.
- Bartoszek K, Sagitov S. 2015a. Phylogenetic confidence intervals for the optimal trait value. *Journal of Applied Probability* 52(04):1115–1132 DOI 10.1239/jap/1450802756.
- Bartoszek K, Sagitov S. 2015b. A consistent estimator of the evolutionary rate. *Journal of Theoretical Biology* 371:69–78 DOI 10.1016/j.jtbi.2015.01.019.
- Bastide P, Solís-Lemus C, Kriebel R, Sparks KW, Ané C. 2018. Phylogenetic comparative methods on phylogenetic networks with reticulations. *Systematic Biology* 67(4):800–820 DOI 10.1093/sysbio/syy005.
- Baurain D, Brinkmann H, Philippe H. 2006. Lack of resolution in the animal phylogeny: closely spaced cladogeneses or undetected systematic errors? *Molecular Biology and Evolution* 24(1):6–9 DOI 10.1093/molbev/msl137.
- Berthelot C, Villar D, Horvath JE, Odom DT, Flicek P. 2018. Complexity and conservation of regulatory landscapes underlie evolutionary resilience of mammalian gene expression. *Nature Ecology & Evolution* 2(1):152–163 DOI 10.1038/s41559-017-0377-2.
- Betancur R, Naylor GJP, Ortí G. 2014. Conserved genes, sampling error, and phylogenomic inference. *Systematic Biology* 63(2):257–262 DOI 10.1093/sysbio/syt073.
- Blischak PD, Chifman J, Wolfe AD, Kubatko LS. 2018. HyDe: A Python Package for Genome-Scale Hybridization Detection. *Systematic Biology* 67:821–829 DOI 10.1093/sysbio/syy023.

- Bleidorn C. 2017.** Sources of error and incongruence in phylogenomic analyses. In: *Phylogenomics*. Cham: Springer International Publishing, 173–193.
- Blom MPK. 2015.** EAPhy: a flexible tool for high-throughput quality filtering of exon-alignments and data processing for phylogenetic methods. Epub ahead of print 5 August 2015. *PLOS Currents* DOI 10.1371/currents.tol.75134257bd389c04bc1d26d42aa9089f
- Boore J. 2006.** The use of genome-level characters for phylogenetic reconstruction. *Trends in Ecology & Evolution* 21(8):439–446 DOI 10.1016/j.tree.2006.05.009.
- Boore JL, Daehler LL, Brown WM. 1999.** Complete sequence, gene arrangement, and genetic code of mitochondrial DNA of the cephalochordate *Branchiostoma floridae* (Amphioxus). *Molecular Biology and Evolution* 16(3):410–418 DOI 10.1093/oxfordjournals.molbev.a026122.
- Brown JW, Wang N, Smith SA. 2017.** The development of scientific consensus: analyzing conflict and concordance among avian phylogenies. *Molecular Phylogenetics and Evolution* 116:69–77 DOI 10.1016/j.ympev.2017.08.002.
- Brumfield RT, Beerli P, Nickerson DA, Edwards SV. 2003.** The utility of single nucleotide polymorphisms in inferences of population history. *Trends in Ecology & Evolution* 18:249–256 DOI 10.1016/s0169-5347(03)00018-1.
- Bryant D, Bouckaert R, Felsenstein J, Rosenberg NA, RoyChoudhury A. 2012.** Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a full coalescent analysis. *Molecular Biology and Evolution* 29(8):1917–1932 DOI 10.1093/molbev/mss086.
- Buerki S, Baker WJ. 2016.** Collections-based research in the genomics era. *Biological Journal of the Linnean Society* 117(1):5–10 DOI 10.1111/bij.12721.
- Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. 2009.** trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* 25(15):1972–1973 DOI 10.1093/bioinformatics/btp348.
- Castresana J. 2000.** Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Molecular Biology and Evolution* 17(4):540–552 DOI 10.1093/oxfordjournals.molbev.a026334.
- Chakrabarty P. 2010.** Genotypes: a concept to help integrate molecular phylogenetics and taxonomy. *Zootaxa* 2632(1):67–68 DOI 10.11646/zootaxa.2632.1.4.
- Chen MY, Liang D, Zhang P. 2015.** Selecting question-specific genes to reduce incongruence in phylogenomics: a case study of jawed vertebrate backbone phylogeny. *Systematic Biology* 64(6):1104–1120 DOI 10.1093/sysbio/syv059.
- Cheng S, Melkonian M, Smith SA, Brockington S, Archibald JM, Delaux P-M, Li F-W, Melkonian B, Mavrodiev EV, Fu WSY, Yang H, Soltis DE, Graham SW, Soltis PS, Liu X, Xu X, Wong GK-S. 2018.** 10KP: A phylodiverse genome sequencing plan. *GigaScience* 7(3):1–9 DOI 10.1093/gigascience/giy013.
- Chifman J, Kubatko L. 2014.** Quartet inference from SNP data under the coalescent model. *Bioinformatics* 30(23):3317–3324 DOI 10.1093/bioinformatics/btu530.
- Churakov G, Sadasivuni MK, Rosenbloom KR, Huchon D, Brosius J, Schmitz J. 2010.** Rodent evolution: back to the root. *Molecular Biology and Evolution* 27(6):1315–1326 DOI 10.1093/molbev/msq019.
- Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. 2016.** GenBank. *Nucleic Acids Research* 44(D1):D67–D72 DOI 10.1093/nar/gkv1276.
- Cloutier A, Sackton TB, Grayson P, Clamp M, Baker AJ, Edwards SV. 2018a.** Whole-genome analyses resolve the phylogeny of flightless birds (Palaeognathae) in the presence of an empirical anomaly zone. *bioRxiv* DOI 10.1101/262949.

- Cloutier A, Sackton TB, Grayson P, Edwards SV, Baker AJ. 2018b. First nuclear genome assembly of an extinct moa species, the little bush moa (*Anomalopteryx didiformis*). *bioRxiv* DOI 10.1101/262816.
- Cohen O, Doron S, Wurtzel O, Dar D, Edelheit S, Karunker I, Mick E, Sorek R. 2016. Comparative transcriptomics across the prokaryotic tree of life. *Nucleic Acids Research* 44(W1):W46–W53 DOI 10.1093/nar/gkw394.
- Costello MJ, Wieczorek J. 2014. Best practice for biodiversity data management and publication. *Biological Conservation* 173:68–73 DOI 10.1016/j.biocon.2013.10.018.
- Crawford FW, Suchard MA. 2013. Diversity, disparity, and evolutionary rate estimation for unresolved Yule trees. *Systematic Biology* 62(3):439–455 DOI 10.1093/sysbio/syt010.
- Cutter AD. 2013. Integrating phylogenetics, phylogeography and population genetics through genomes and evolutionary theory. *Molecular Phylogenetics and Evolution* 69(3):1172–1185 DOI 10.1016/j.ympev.2013.06.006.
- Dalquen DA, Zhu T, Yang Z. 2017. Maximum likelihood implementation of an isolation-with-migration model for three species. *Systematic Biology* 66(3):379–398 DOI 10.1093/sysbio/syw063.
- Darriba D, Flouri T, Stamatakis A. 2018. The state of software for evolutionary biology. *Molecular Biology and Evolution* 35(5):1037–1046 DOI 10.1093/molbev/msy014.
- Davies TJ, Kraft NJB, Salamin N, Wolkovich EM. 2012. Incompletely resolved phylogenetic trees inflate estimates of phylogenetic conservatism. *Ecology* 93(2):242–247 DOI 10.1890/11-1360.1.
- Dayrat B. 2005. Towards integrative taxonomy. *Biological Journal of the Linnean Society* 85(3):407–415 DOI 10.1111/j.1095-8312.2005.00503.x.
- DeBiasse MB, Ryan JF. 2018. Phylotocol: Promoting transparency and overcoming bias in phylogenetics. *PeerJ Preprints* 6:e26585v4 DOI 10.7287/peerj.preprints.26585v4.
- De Maio N, Schlötterer C, Kosiol C. 2013. Linking great apes genome evolution across time scales using polymorphism-aware phylogenetic models. *Molecular Biology and Evolution* 30(10):2249–2262 DOI 10.1093/molbev/mst131.
- de Queiroz A, Gatesy J. 2007. The supermatrix approach to systematics. *Trends in Ecology & Evolution* 22(1):34–41 DOI 10.1016/j.tree.2006.10.002.
- de Vienne DM, Ollier S, Aguilera G. 2012. Phylo-MCOA: a fast and efficient method to detect outlier genes and species in phylogenomics using multiple co-inertia analysis. *Molecular Biology and Evolution* 29(6):1587–1598 DOI 10.1093/molbev/msr317.
- DeBolt S. 2010. Copy number variation shapes genome diversity in arabidopsis over immediate family generational scales. *Genome Biology and Evolution* 2(10):441–453 DOI 10.1093/gbe/evq033.
- Degnan JH. 2018. Modeling hybridization under the network multispecies coalescent. *Systematic Biology* 67(5):786–799 DOI 10.1093/sysbio/syy040.
- Degnan JH, Rosenberg NA. 2009. Gene tree discordance, phylogenetic inference and the multispecies coalescent. *Trends in Ecology & Evolution* 24(6):332–340 DOI 10.1016/j.tree.2009.01.009.
- Dell Ampio E, Meusemann K, Szucsich NU, Peters RS, Meyer B, Borner J, Petersen M, Aberer AJ, Stamatakis A, Walz MG, Minh BQ, Haeseler von A, Ebersberger I, Pass GN, Misof B. 2013. Decisive data sets in phylogenomics: lessons from studies on the phylogenetic relationships of primarily wingless insects. *Molecular Biology and Evolution* 31(1):239–249 DOI 10.1093/molbev/mst196.

- Delsuc F, Brinkmann H, Philippe H. 2005. Phylogenomics and the reconstruction of the tree of life. *Nature Reviews Genetics* 6(5):361–375 DOI 10.1038/nrg1603.
- dos Reis M, Donoghue PCJ, Yang Z. 2016. Bayesian molecular clock dating of species divergences in the genomics era. *Nature Reviews Genetics* 17(2):71–80 DOI 10.1038/nrg.2015.8.
- Douzery EJP, Scornavacca C, Romiguier J, Belkhir K, Galtier N, Delsuc F, Ranwez V. 2014. OrthoMaM v8: a database of orthologous exons and coding sequences for comparative genomics in mammals. *Molecular Biology and Evolution* 31(7):1923–1928 DOI 10.1093/molbev/msu132.
- Doyle JJ. 1992. Gene trees and species trees: molecular systematics as one-character taxonomy. *Systematic Botany* 17(1):144–163 DOI 10.2307/2419070.
- Drew BT. 2013. Data deposition: missing data mean holes in tree of life. *Nature* 493(7432):305–305 DOI 10.1038/493305f.
- Drew BT, Gazis R, Cabezas P, Swithers KS, Deng J, Rodriguez R, Katz LA, Crandall KA, Hibbett DS, Soltis DE. 2013. Lost branches on the tree of life. *PLOS Biology* 11(9):e1001636 DOI 10.1371/journal.pbio.1001636.
- Dunn CW, Howinson M, Zapata F. 2013. Agalma: an automated phylogenomics workflow. *BMC Bioinformatics* 14(1):330 DOI 10.1186/1471-2105-14-330.
- Edwards SV. 2009a. Natural selection and phylogenetic analysis. *Proceedings of the National Academy of Sciences of the United States of America* 106(22):8799–8800 DOI 10.1073/pnas.0904103106.
- Edwards SV. 2009b. Is a new and general theory of molecular systematics emerging? *Evolution* 63(1):1–19 DOI 10.1111/j.1558-5646.2008.00549.x.
- Edwards SV, Beerli P. 2000. Perspective: gene divergence, population divergence, and the variance in coalescence time in phylogeographic studies. *Evolution* 54(6):1839–1854 DOI 10.1111/j.0014-3820.2000.tb01231.x.
- Edwards SV, Cloutier A, Baker AJ. 2017. Conserved nonexonic elements: a novel class of marker for phylogenomics. *Systematic Biology* 66(6):1028–1044 DOI 10.1093/sysbio/syx058.
- Edwards SV, Liu L, Pearl DK. 2007. High-resolution species trees without concatenation. *Proceedings of the National Academy of Sciences of the United States of America* 104(14):5936–5941 DOI 10.1073/pnas.0607004104.
- Edwards SV, Potter S, Schmitt CJ, Bragg JG, Moritz C. 2016a. Reticulation, divergence, and the phylogeography–phylogenetics continuum. *Proceedings of the National Academy of Sciences of the United States of America* 113(29):8025–8032 DOI 10.1073/pnas.1601066113.
- Edwards SV, Xi Z, Janke A, Faircloth BC, McCormack JE, Glenn TC, Zhong B, Wu S, Lemmon EM, Lemmon AR, Leaché AD, Liu L, Davis CC. 2016b. Implementing and testing the multispecies coalescent model: a valuable paradigm for phylogenomics. *Molecular Phylogenetics and Evolution* 94:447–462 DOI 10.1016/j.ympev.2015.10.027.
- Ellegren N, Galtier N. 2016. Determinants of genetic diversity. *Nature Reviews Genetics* 17(7):422–433 DOI 10.1038/nrg.2016.58.
- Elmer KR, Meyer A. 2011. Adaptation in the age of ecological genomics: insights from parallelism and convergence. *Trends in Ecology & Evolution* 26(6):298–306 DOI 10.1016/j.tree.2011.02.008.
- Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust demographic inference from genomic and SNP data. *PLOS Genetics* 9(10):e1003905 DOI 10.1371/journal.pgen.1003905.
- Faircloth BC, McCormack JE, Crawford NG, Harvey MG, Brumfield RT, Glenn TC. 2012. Ultraconserved elements anchor thousands of genetic markers spanning multiple evolutionary timescales. *Systematic Biology* 61(5):717–726 DOI 10.1093/sysbio/sys004.

- Faurby S, Svenning JC. 2015. A species-level phylogeny of all extant and late Quaternary extinct mammals using a novel heuristic-hierarchical Bayesian approach. *Molecular Phylogenetics and Evolution* 84:14–26 DOI 10.1016/j.ympev.2014.11.001.
- Felsenstein J. 1973. Maximum-likelihood estimation of evolutionary trees from continuous characters. *American Journal of Human Genetics* 25:471–492.
- Felsenstein J. 1985. Phylogenies and the comparative method. *American Naturalist* 125(1):1–15 DOI 10.1086/284325.
- Felsenstein J. 1988. Phylogenies from molecular sequences: inference and reliability. *Annual Review of Genetics* 22(1):521–565 DOI 10.1146/annurev.ge.22.120188.002513.
- Felsenstein J. 2012. A comparative method for both discrete and continuous characters using the threshold model. *American Naturalist* 179(2):145–156 DOI 10.1086/663681.
- Fernández R, Laumer CE, Vahtera V, Libro S, Kaluziak S, Sharma PP, Pérez-Morro AR, Edgecombe GD, Giribert G. 2014. Evaluating topological conflict in centipede phylogeny using transcriptomic data sets. *Molecular Biology and Evolution* 31(6):1500–1513 DOI 10.1093/molbev/msu108.
- Figuet E, Ballenghien M, Romiguier J, Galtier N. 2015. Biased gene conversion and GC-content evolution in the coding sequences of reptiles and vertebrates. *Genome Biology and Evolution* 7(1):240–250 DOI 10.1093/gbe/evu277.
- Fitch WM. 1970. Distinguishing homologous from analogous proteins. *Systematic Zoology* 19(2):99–113 DOI 10.2307/2412448.
- Fong JJ, Brown JM, Fujita MK, Boussau B. 2012. A Phylogenomic approach to vertebrate phylogeny supports a turtle-archosaur affinity and a possible paraphyletic lissamphibia. *PLOS ONE* 7(11):e48990 DOI 10.1371/journal.pone.0048990.
- Frantzeskakis L, Kracher B, Kusch S, Yoshikawa-Maekawa M, Bauer S, Pedersen C, Spanu PD, Maekawa T, Schulze-Lefert P, Panstruga R. 2018. Signatures of host specialization and a recent transposable element burst in the dynamic one-speed genome of the fungal barley powdery mildew pathogen. *BMC Genomics* 19(1):381 DOI 10.1186/s12864-018-4750-6.
- Fredman D, White SJ, Potter S, Eichler EE, Dunnen JTD, Brookes AJ. 2004. Complex SNP-related sequence variation in segmental genome duplications. *Nature Genetics* 36(8):861–866 DOI 10.1038/ng1401.
- Garrick RC, Bonatelli IAS, Hyseni C, Morales A, Pelletier TA, Perez MF, Rice E, Satler JD, Symula RE, Thomé MTC, Carstens BC. 2015. The evolution of phylogeographic data sets. *Molecular Ecology* 24(6):1164–1171 DOI 10.1111/mec.13108.
- Garrick RC, Sunnucks P, Dyer RJ. 2010. Nuclear gene phylogeography using PHASE: dealing with unresolved genotypes, lost alleles, and systematic bias in parameter estimation. *BMC Evolutionary Biology* 10(1):118 DOI 10.1186/1471-2148-10-118.
- Genome 10K Community of Scientists. 2009. A proposal to obtain whole-genome sequence for 10,000 vertebrate species. *Journal of Heredity* 100(6):659–674 DOI 10.1093/jhered/esp086.
- Gernhard T. 2008a. The conditioned reconstructed process. *Journal of Theoretical Biology* 253(4):769–778 DOI 10.1016/j.jtbi.2008.04.005.
- Gernhard T. 2008b. New analytic results for speciation times in neutral models. *Bulletin of Mathematical Biology* 70(4):1082–1097 DOI 10.1007/s11538-007-9291-0.
- Ghiurcuta CG, Moret BME. 2014. Evaluating synteny for improved comparative studies. *Bioinformatics* 30(12):i9–i18 DOI 10.1093/bioinformatics/btu259.
- Gogarten JP, Townsend JP. 2005. Horizontal gene transfer, genome innovation and evolution. *Nature Reviews Microbiology* 3:679–687 DOI 10.1038/nrmicro1204.

- Goodwin S, McPherson JD, McCombie WR. 2016.** Coming of age: ten years of next-generation sequencing technologies. *Nature Reviews Genetics* **17**(6):333–351 DOI [10.1038/nrg.2016.49](https://doi.org/10.1038/nrg.2016.49).
- Gratton P, Marta S, Bocksberger G, Winter M, Trucchi E, Kühl H. 2017.** A world of sequences: can we use georeferenced nucleotide databases for a robust automated phylogeography? *Journal of Biogeography* **44**(2):475–486 DOI [10.1111/jbi.12786](https://doi.org/10.1111/jbi.12786).
- Graybeal A. 1998.** Is it better to add taxa or characters to a difficult phylogenetic problem? *Systematic Biology* **47**(1):9–17 DOI [10.1080/106351598260996](https://doi.org/10.1080/106351598260996).
- Griffiths RC, Marjoram P. 1996.** An ancestral recombination graph. In: Donnelly P, Tavaré S, eds. *IMA Volume on Mathematical Population Genetics*. New York: Springer-Verlag, 257–270.
- Grigoriev IV, Nikitin R, Haridas S, Kuo A, Ohm R, Otilar R, Riley R, Salamov A, Zhao X, Korzeniewski F, Smirnova T, Nordberg H, Dubchak I, Shabalov I. 2014.** MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic Acids Research* **42**(D1):D699–D704 DOI [10.1093/nar/gkt1183](https://doi.org/10.1093/nar/gkt1183).
- Gronau I, Hubisz MJ, Gulko B, Danko CG, Siepel A. 2011.** Bayesian inference of ancient human demography from individual genome sequences. *Nature Genetics* **43**:1031–1034 DOI [10.1038/ng.937](https://doi.org/10.1038/ng.937).
- Guang A, Zapata F, Howison M, Lawrence CE, Dunn CW. 2016.** An integrated perspective on phylogenetic workflows. *Trends in Ecology and Evolution* **31**(2):116–126 DOI [10.1016/j.tree.2015.12.007](https://doi.org/10.1016/j.tree.2015.12.007).
- Guerrero RF, Hahn MW. 2018.** Quantifying the risk of hemiplasy in phylogenetic inference. *Proceedings of the National Academy of Sciences of the United States of America* **115**(50):12787–12792 DOI [10.1073/pnas.1811268115](https://doi.org/10.1073/pnas.1811268115).
- Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010.** New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Systematic Biology* **59**(3):307–321 DOI [10.1093/sysbio/syq010](https://doi.org/10.1093/sysbio/syq010).
- Gusfield D. 2015.** *Persistent phylogeny*. New York: ACM Press, 443–451.
- Hahn MW. 2018.** *Molecular population genomics*. Sunderland: Oxford University Press.
- Hahn MW, Nakhleh L. 2016.** Irrational exuberance for resolved species trees. *Evolution* **70**(1):7–17 DOI [10.1111/evo.12832](https://doi.org/10.1111/evo.12832).
- Hammer MF, Woerner AE, Mendez FL, Watkins JC, Wall JD. 2011.** Genetic evidence for archaic admixture in Africa. *Proceedings of the National Academy of Sciences of the United States of America* **108**(37):15123–15128 DOI [10.1073/pnas.1109300108](https://doi.org/10.1073/pnas.1109300108).
- Hansen TF. 1997.** Stabilizing selection and the comparative analysis of adaptation. *Evolution* **51**(5):1341–1351 DOI [10.1111/j.1558-5646.1997.tb01457.x](https://doi.org/10.1111/j.1558-5646.1997.tb01457.x).
- Harmon LJ, Baumes J, Hughes C, Soberón J, Specht CD, Tumer W, Lisle C, Thacker RW. 2013.** Arbor: comparative analysis workflows for the tree of life. *PLOS Currents* **5** DOI [10.1371/currents.tol.099161de5eabdee073fd3d21a44518dc](https://doi.org/10.1371/currents.tol.099161de5eabdee073fd3d21a44518dc).
- Harvey MG, Smith BT, Glenn TC, Faircloth BC, Brumfield RT. 2016.** Sequence capture versus restriction site associated DNA sequencing for shallow systematics. *Systematic Biology* **65**(5):910–924 DOI [10.1093/sysbio/syw036](https://doi.org/10.1093/sysbio/syw036).
- He D, Sierra R, Pawlowski J, Baldauf SL. 2016.** Reducing long-branch effects in multi-protein data uncovers a close relationship between *Alveolata* and *Rhizaria*. *Molecular Phylogenetics and Evolution* **101**:1–7 DOI [10.1016/j.ympev.2016.04.033](https://doi.org/10.1016/j.ympev.2016.04.033).
- Heath TA, Huelsenbeck JP, Stadler T. 2014.** The fossilized birth-death process for coherent calibration of divergence-time estimates. *Proceedings of the National Academy of Sciences of the United States of America* **111**(29):E2957–E2966 DOI [10.1073/pnas.1319091111](https://doi.org/10.1073/pnas.1319091111).

- Hedtke SM, Townsend TM, Hillis DM. 2006. Resolution of phylogenetic conflict in large data sets by increased taxon sampling. *Systematic Biology* 55(3):522–529 DOI 10.1080/10635150600697358.
- Heled J, Drummond AJ. 2010. Bayesian inference of species trees from multilocus data. *Molecular Biology and Evolution* 27(3):570–580 DOI 10.1093/molbev/msp274.
- Hey J. 1992. Using phylogenetic trees to study speciation and extinction. *Evolution* 46(3):627–640 DOI 10.1111/j.1558-5646.1992.tb02071.x.
- Hey J. 2010. Isolation with migration models for more than two populations. *Molecular Biology and Evolution* 27(4):905–920 DOI 10.1093/molbev/msp296.
- Hey J, Nielsen R. 2004. Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*. *Genetics* 167(2):747–760 DOI 10.1534/genetics.103.024182.
- Hey J, Nielsen R. 2007. Integration within the Felsenstein equation for improved Markov chain Monte Carlo methods in population genetics. *Proceedings of the National Academy of Sciences of the United States of America* 104(8):2785–2790 DOI 10.1073/pnas.0611164104.
- Hibbins MS, Hahn MW. 2018. The timing and direction of introgression under the multispecies network coalescent. *bioRxiv* DOI 10.1101/328575.
- Hiller M, Schaar BT, Indjeian VB, Kingsley DM, Hagey LR, Bejerano G. 2012. A “forward genomics” approach links genotype to phenotype using independent phenotypic losses among related species. *Cell Reports* 2(4):817–823 DOI 10.1016/j.celrep.2012.08.032.
- Hillis DM. 1996. Inferring complex phylogenies. *Nature* 383(6596):130–131 DOI 10.1038/383130a0.
- Hillis DM. 1998. Taxonomic sampling, phylogenetic accuracy, and investigator bias. *Systematic Biology* 47(1):3–8 DOI 10.1080/106351598260987.
- Hinchliff CE, Smith SA, Allman JF, Burleigh JG, Chaudhary R, Coghill LM, Crandall KA, Deng J, Drew BT, Gazis R, Gude K, Hibbett DS, Katz LA, Laughinghouse HDIV, McTavish EJ, Midford PE, Owen CL, Ree RH, Rees JA, Soltis DE, Williams T, Cranston KA. 2015. Synthesis of phylogeny and taxonomy into a comprehensive tree of life. *Proceedings of the National Academy of Sciences of the United States of America* 112(41):12764–12769 DOI 10.1073/pnas.1423041112.
- Hoban S, Kelley JL, Lotterhos KE, Antolin MF, Bradburd G, Lowry DB, Poss ML, Reed LK, Storfer A, Whitlock MC. 2016. Finding the genomic basis of local adaptation: pitfalls, practical solutions, and future directions. *American Naturalist* 188(4):379–397 DOI 10.1086/688018.
- Hobolth A, Christensen OF, Mailund T, Schierup MH. 2007. Genomic relationships and speciation times of human, chimpanzee, and gorilla inferred from a coalescent hidden Markov model. *PLOS Genetics* 3(2):e7 DOI 10.1371/journal.pgen.0030007.
- Huang HT, He QI, Kubatko LS, Knowles LL. 2010. Sources of error inherent in species-tree estimation: impact of mutational and coalescent effects on accuracy and implications for choosing among different methods. *Systematic Biology* 59(5):573–583 DOI 10.1093/sysbio/syq047.
- Huang H, Sukumaran J, Smith SA, Knowles LL. 2017. Cause of gene tree discord? Distinguishing incomplete lineage sorting and lateral gene transfer in phylogenetics. *PeerJ Preprints* 5:3489v1 DOI 10.7287/peerj.preprints.3489v1.
- Huang W, Zhou G, Marchand M, Ash JR, Morris D, Van Dooren P, Brown JM, Gallivan KA, Wilgenbusch JC. 2016. TreeScaper: visualizing and extracting phylogenetic signal from sets of trees. *Molecular Biology and Evolution* 33(12):3314–3316 DOI 10.1093/molbev/msw196.

- Huber KT, Oxelman B, Lott M, Moulton V. 2006. Reconstructing the evolutionary history of polyploids from multilabeled trees. *Molecular Biology and Evolution* 23(9):1784–1791 DOI 10.1093/molbev/msl045.
- Hunt G. 2013. Testing the link between phenotypic evolution and speciation: an integrated palaeontological and phylogenetic analysis. *Methods in Ecology and Evolution* 4(8):714–723 DOI 10.1111/2041-210X.12085.
- Hunt G, Slater G. 2016. Integrating paleontological and phylogenetic approaches to macroevolution. *Annual Review of Ecology, Evolution, and Systematics* 47(1):189–213 DOI 10.1146/annurev-ecolsys-112414-054207.
- Huson DH, Bryant D. 2006. Application of phylogenetic networks in evolutionary studies. *Molecular Biology and Evolution* 23(2):254–267 DOI 10.1093/molbev/msj030.
- Hykin SM, Bi K, McGuire JA. 2015. Fixing formalin: a method to recover genomic-scale DNA sequence data from formalin-fixed museum specimens using high-throughput sequencing. *PLOS ONE* 10(10):e0141579 DOI 10.1371/journal.pone.0141579.
- Iqbal Z, Caccamo M, Turner I, Flicek P, McVean G. 2012. De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nature Genetics* 44(2):226–232 DOI 10.1038/ng.1028.
- Irisarri I, Baurain D, Brinkmann H, Delsuc F, Sire J-Y, Kupfer A, Petersen J, Jarek M, Meyer A, Vences M, Philippe H. 2017. Phylotranscriptomic consolidation of the jawed vertebrate timetree. *Nature Ecology & Evolution* 1(9):1370–1378 DOI 10.1038/s41559-017-0240-5.
- Jackson ND, Morales AE, Carstens BC, O'Meara BC. 2017. PHRAPL: Phylogeographic Inference Using Approximate Likelihoods. *Systematic Biology* 66:1045–1053 DOI 10.1093/sysbio/syx001.
- Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SYW, Faircloth BC, Nabholz B, Howard JT, Suh A, Weber CC, da Fonseca RR, Li J, Zhang F, Li H, Zhou L, Narula N, Liu L, Ganapathy G, Boussau B, Bayzid MS, Zavidovych V, Subramanian S, Gabaldon T, Capella-Gutierrez S, Huerta-Cepas J, Rekepalli B, Munch K, Schierup M, Lindow B, Warren WC, Ray D, Green RE, Bruford MW, Zhan X, Dixon A, Li S, Li N, Huang Y, Derryberry EP, Bertelsen MF, Sheldon FH, Brumfield RT, Mello CV, Lovell PV, Wirthlin M, Schneider MPC, Prosdocimi F, Samaniego JA, Velazquez AMV, Alfaro-Nunez A, Campos PF, Petersen B, Sicheritz-Ponten T, Pas A, Bailey T, Scofield P, Bunce M, Lambert DM, Zhou Q, Perelman P, Driskell AC, Shapiro B, Xiong Z, Zeng Y, Liu S, Li Z, Liu B, Wu K, Xiao J, Yinqi X, Zheng Q, Zhang Y, Yang H, Wang J, Smeds L, Rheindt FE, Braun M, Fjeldsa J, Orlando L, Barker FK, Jonsson KA, Johnson W, Koepfli K-P, O'Brien S, Haussler D, Ryder OA, Rahbek C, Willerslev E, Graves GR, Glenn TC, McCormack J, Burt D, Ellegren H, Alstrom P, Edwards SV, Stamatakis A, Mindell DP, Cracraft J, Braun EL, Warnow T, Jun W, Gilbert MTP, Zhang G. 2014. Whole-genome analyses resolve early branches in the tree of life of modern birds. *Science* 346(6215):1320–1331 DOI 10.1126/science.1253451.
- Jeffroy O, Brinkmann H, Delsuc F, Philippe H. 2006. Phylogenomics: the beginning of incongruence? *Trends in Genetics* 22(4):225–231 DOI 10.1016/j.tig.2006.02.003.
- Jetz W, Thomas GH, Joy JB, Hartmann K, Mooers AO. 2012. The global diversity of birds in space and time. *Nature* 491(7424):444–448 DOI 10.1038/nature11631.
- Jhwueng DC, O'Meara B. 2015. Trait evolution on phylogenetic networks. *BioRxiv* DOI 10.1101/023986.
- Johnson SS, Zaikova E, Goerlitz DS, Bai Y, Tighe SW. 2017. Real-time DNA sequencing in the antarctic dry valleys using the oxford nanopore sequencer. *Journal of Biomolecular Techniques* 28(1):2–7 DOI 10.7171/jbt.17-2801-009.

- Jombart T, Kendall M, Almagro-Garcia J, Colijn C. 2017. TREESPACE: Statistical exploration of landscapes of phylogenetic trees. *Molecular Ecology Resources* 17(6):1385–1392 DOI 10.1111/1755-0998.12676.
- Jones GR. 2017. Algorithmic improvements to species delimitation and phylogeny estimation under the multispecies coalescent. *Journal of Mathematical Biology* 74(1–2):447–467 DOI 10.1007/s00285-016-1034-0.
- Jones GR. 2018. Divergence estimation in the presence of incomplete lineage sorting and migration. *Systematic Biology* 68(1):19–31 DOI 10.1093/sysbio/syy041.
- Jones MR, Good JM. 2016. Targeted capture in evolutionary and ecological genomics. *Molecular Ecology* 25(1):185–202 DOI 10.1111/mec.13304.
- Jones GR, Aydin Z, Oxelman B. 2015. DISSECT: an assignment-free Bayesian discovery method for species delimitation under the multispecies coalescent. *Bioinformatics* 31(7):991–998 DOI 10.1093/bioinformatics/btu770.
- Jones GR, Sagitov S, Oxelman B. 2013. Statistical inference of allopolyploid species networks in the presence of incomplete lineage sorting. *Systematic Biology* 62(3):467–478 DOI 10.1093/sysbio/syt012.
- Kainer D, Lanfear R. 2015. The effects of partitioning on phylogenetic inference. *Molecular Biology and Evolution* 32(6):1611–1627 DOI 10.1093/molbev/msv026.
- Kaiser VB, van Tuinen M, Ellegren H. 2007. Insertion events of CR1 retrotransposable elements elucidate the phylogenetic branching order in galliform birds. *Molecular Biology and Evolution* 24(1):338–347 DOI 10.1093/molbev/msl164.
- Khost DE, Eickbush DG, Larracuente AM. 2017. Single-molecule sequencing resolves the detailed structure of complex satellite DNA loci in *Drosophila melanogaster*. *Genome Research* 27(5):709–721 DOI 10.1101/gr.213512.116.
- Kidd DM. 2010. Geophylogenies and the map of life. *Systematic Biology* 59(6):741–752 DOI 10.1093/sysbio/syq043.
- Kim J. 1996. General inconsistency conditions for maximum parsimony: effects of branch lengths and increasing numbers of taxa. *Systematic Biology* 46(3):363–374 DOI 10.1093/sysbio/45.3.363.
- King N, Rokas A. 2017. Embracing uncertainty in reconstructing early animal evolution. *Current Biology* 27(19):R1081–R1088 DOI 10.1016/j.cub.2017.08.054.
- Klopfstein S, Massingham T, Goldman N. 2017. More on the best evolutionary rate for phylogenetic analysis. *Systematic Biology* 66(5):769–785 DOI 10.1093/sysbio/syx051.
- Knowles LL, Huang H, Sukumaran J, Smith SA. 2018. A matter of phylogenetic scale: distinguishing incomplete lineage sorting from lateral gene transfer as the cause of gene tree discord in recent versus deep diversification histories. *American Journal of Botany* 105(3):376–384 DOI 10.1002/ajb2.1064.
- Kocot KM, Struck TH, Merkel J, Waits DS, Todt C, Brannock PM, Weese DA, Cannon JT, Moroz LL, Lieb B, Halanych KM. 2017. Phylogenomics of lophotrochozoa with consideration of systematic error. *Systematic Biology* 66(2):256–282 DOI 10.1093/sysbio/syw079.
- Kowada LAB, Doerr D, Dantas S, Stoye J. 2016. New genome similarity measures based on conserved gene adjacencies. In: Singh M, ed. *Research in Computational Molecular Biology. RECOMB 2016. Lecture Notes in Computer Science*. Vol. 9649. Cham: Springer.
- Kozlov AM, Aberer AJ, Stamatakis A. 2015. ExaML version 3: a tool for phylogenomic analyses on supercomputers. *Bioinformatics* 31(15):2577–2579 DOI 10.1093/bioinformatics/btv184.

- Kriegs JO, Zemann A, Churakov G, Matzke A, Ohme M, Zischler H, Brosius J, Kryger U, Schmitz J. 2010. Retroposon insertions provide insights into deep lagomorph evolution. *Molecular Biology and Evolution* 27(12):2678–2681 DOI 10.1093/molbev/msq162.
- Ksepka DT, Benton MJ, Carrano MT, Gandolfo MA, Head JJ, Hermesen EJ, Joyce WG, Lamm KS, Patané JSL, Phillips MJ, Polly PD, van Tuinen M, Ware JL, Warnock RCM, Parham JF. 2011. Synthesizing and databasing fossil calibrations: divergence dating and beyond. *Biology Letters* 7(6):801–803 DOI 10.1098/rsbl.2011.0356.
- Kubatko LS, Degnan JH. 2007. Inconsistency of phylogenetic estimates from concatenated data under coalescence. *Systematic Biology* 56(1):17–24 DOI 10.1080/10635150601146041.
- Kubatko LS, Gibbs HL, Bloomquist EW. 2011. Inferring species-level phylogenies and taxonomic distinctiveness using multilocus data in *Sistrurus rattlesnakes*. *Systematic Biology* 60(4):393–409 DOI 10.1093/sysbio/syr011.
- Kuhn TS, Mooers AØ, Thomas GH. 2011. A simple polytomy resolver for dated phylogenies. *Methods in Ecology and Evolution* 2(5):427–436 DOI 10.1111/j.2041-210X.2011.00103.x.
- Kumar S, Filipski AJ, Battistuzzi FU, Kosakovsky Pond SL, Tamura K. 2012. Statistics and truth in phylogenomics. *Molecular Biology and Evolution* 29(2):457–472 DOI 10.1093/molbev/msr202.
- Kumar S, Jones M, Koutsovoulos G, Clarke M, Blaxter M. 2013. Blobology: exploring raw genome data for contaminants, symbionts and parasites using taxon-annotated GC-coverage plots. *Frontiers in Genetics* 4:237 DOI 10.3389/fgene.2013.00237.
- Kunin V, Goldovsky L, Darzentas N. 2005. The net of life: reconstructing the microbial phylogenetic network. *Genome Research* 15(7):954–959 DOI 10.1101/gr.3666505.
- Lamichhaney S, Berglund B, Almén MS, Maqbool K, Grabherr M, Martinez-Barrio A, Promerova M, Rubin CJ, Wang C, Zamani N, Grant BR, Grant PR, Webster MT, Andersson L. 2015. Evolution of Darwin's finches and their beaks revealed by genome sequencing. *Nature* 518(7539):371–375 DOI 10.1038/nature14181.
- Lanier HC, Knowles LL. 2012. Is recombination a problem for species-tree analyses? *Systematic Biology* 61(4):691–701 DOI 10.1093/sysbio/syr128.
- Lanier HC, Knowles LL. 2015. Applying species-tree analyses to deep phylogenetic histories: challenges and potential suggested from a survey of empirical phylogenetic studies. *Molecular Phylogenetics and Evolution* 83:191–199 DOI 10.1016/j.ympev.2014.10.022.
- Lanfear R, Calcott B, Kainer D, Mayer C, Stamatakis A. 2014. Selecting optimal partitioning schemes for phylogenomic datasets. *BMC Evolutionary Biology* 14(1):82 DOI 10.1186/1471-2148-14-82.
- Larson-Johnson K. 2016. Phylogenetic investigation of the complex evolutionary history of dispersal mode and diversification rates across living and fossil Fagales. *New Phytologist* 209(1):418–435 DOI 10.1111/nph.13570.
- Lartillot N, Rodrigue N, Stubbs D, Richer J. 2013. PhyloBayes MPI: phylogenetic reconstruction with infinite mixtures of profiles in a parallel environment. *Systematic Biology* 62(4):611–615 DOI 10.1093/sysbio/syt022.
- Leaché AD, Chavez AS, Jones LN, Grummer JA, Gottscho AD, Linkem CW. 2015. Phylogenomics of phrynosomatid lizards: conflicting signals from sequence capture versus restriction site associated DNA sequencing. *Genome Biology and Evolution* 7(3):706–719 DOI 10.1093/gbe/evv026.
- Leaché AD, Harris RB, Rannala B, Yanz Z. 2014. The influence of gene flow on species tree estimation: a simulation study. *Systematic Biology* 63(1):17–30 DOI 10.1093/sysbio/syt049.

- Leaché AD, Oaks JR. 2017. The utility of single nucleotide polymorphism (SNP) data in phylogenetics. *Annual Review of Ecology, Evolution, and Systematics* 48(1):69–84 DOI 10.1146/annurev-ecolsys-110316-022645.
- Leigh JW, Lapointe F-J, Lopez P, Baptiste E. 2011. Evaluating phylogenetic congruence in the post-genomic era. *Genome Biology and Evolution* 3:571–587 DOI 10.1093/gbe/evr050.
- Lelieveld SH, Spielman M, Mundlos S, Veltman JA, Gilissen C. 2015. Comparison of exome and genome sequencing technologies for the complete capture of protein-coding regions. *Human Mutation* 36(8):815–822 DOI 10.1002/humu.22813.
- Lemmon AR, Emme SA, Lemmon EM. 2012. Anchored hybrid enrichment for massively high-throughput phylogenomics. *Systematic Biology* 61(5):727–744 DOI 10.1093/sysbio/sys049.
- Lemmon EM, Lemmon AR. 2013. High-throughput genomic data in systematics and phylogenetics. *Annual Review of Ecology, Evolution, and Systematics* 44(1):99–121 DOI 10.1146/annurev-ecolsys-110512-135822.
- Leonardi M, Librado P, Der Sarkissian C, Schubert M, Alfarhan AH, Alquraishi SA, Al-Rasheid KAS, Gamba C, Willerslev E, Orlando L. 2017. Evolutionary patterns and processes: lessons from ancient DNA. *Systematic Biology* 66:e1–e29.
- Leprevost FV, Barbosa VC, Francisco EL, Perez-Riverol Y, Carvalho PC. 2014. On best practices in the development of bioinformatics software. *Frontiers in Genetics* 5:199 DOI 10.3389/fgene.2014.00199.
- Li G, Figueiró HV, Eizirik E, Murphy WJ. 2018. Recombination-aware phylogenomics unravels the complex divergence of hybridizing species. *bioRxiv* DOI 10.1101/485904.
- Liang Y, Liao Bo, Zhu W. 2017. An improved binary differential evolution algorithm to infer tumor phylogenetic trees. *BioMed Research International* 2017:13.
- Lin Y, Hu F, Tang J, Moret BM. 2013. Maximum likelihood phylogenetic reconstruction from high-resolution whole-genome data and a tree of 68 eukaryotes. In: *Pacific Symposium on Biocomputing*, Singapore: World Scientific, 285–296.
- Lischer HE, Excoffier L, Heckel G. 2014. Ignoring heterozygous sites biases phylogenomic estimates of divergence times: implications for the evolutionary history of *Microtus* voles. *Molecular Biology and Evolution* 31(4):817–831 DOI 10.1093/molbev/mst271.
- Liu L, Pearl DK. 2007. Species trees from gene trees: reconstructing Bayesian posterior distributions of a species phylogeny using estimated gene tree distributions. *Systematic Biology* 56(3):504–514 DOI 10.1080/10635150701429982.
- Liu L, Wu S, Yu L. 2015. Coalescent methods for estimating species trees from phylogenomic data. *Journal of Systematics and Evolution* 53(5):380–390 DOI 10.1111/jse.12160.
- Liu L, Xi Z, Davis CC. 2015. Coalescent methods are robust to the simultaneous effects of long branches and incomplete lineage sorting. *Molecular Biology and Evolution* 32(3):791–805 DOI 10.1093/molbev/msu331.
- Liu L, Xi Z, Wu S, Davis CC, Edwards SV. 2015. Estimating phylogenetic trees from genome-scale data. *Annals of the New York Academy of Sciences* 1360(1):36–53 DOI 10.1111/nyas.12747.
- Liu L, Yu L, Edwards SV. 2010. A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. *BMC Evolutionary Biology* 10(1):302 DOI 10.1186/1471-2148-10-302.
- Liu L, Yu L, Pearl DK, Edwards SV. 2009. Estimating species phylogenies using coalescence times among sequences. *Systematic Biology* 58(5):468–477 DOI 10.1093/sysbio/syp031.

- Liu L, Zhang J, Rheindt FE, Lei F, Qu Y, Wang Y, Sullivan C, Nie W, Wang J, Yang F, Chen J, Edwards SV, Meng J, Wu S. 2017. Genomic evidence reveals a radiation of placental mammals uninterrupted by the KPg boundary. *Proceedings of the National Academy of Science of the United States of America* 114(35):E7282–E7290 DOI 10.1073/pnas.1616744114.
- Long JC. 1991. The genetic structure of admixed populations. *Genetics* 127:417–418.
- Lott M, Spillner A, Huber KT, Moulton V. 2009. PADRE: a package for analyzing and displaying reticulate evolution. *Bioinformatics* 25(9):1199–1200 DOI 10.1093/bioinformatics/btp133.
- Lu H, Giordano F, Ning Z. 2016. Oxford nanopore MinION sequencing and genome assembly. *Genomics, Proteomics and Bioinformatics* 14(5):265–279 DOI 10.1016/j.gpb.2016.05.004.
- Lynch VJ, Bedoya-Reina OC, Ratan A, Sulak M, Drautz-Moses DI, Perry GH, Miller W, Schuster SC. 2015. Elephantid genomes reveal the molecular bases of woolly mammoth adaptations to the arctic. *Cell Reports* 12(2):217–228 DOI 10.1016/j.celrep.2015.06.027.
- Maddison WP. 1997. Gene trees in species trees. *Systematic Biology* 46(3):523–536 DOI 10.1093/sysbio/46.3.523.
- Magee AF, May MR, Moore BR. 2014. The dawn of open access to phylogenetic data. *PLOS ONE* 9(10):e110268 DOI 10.1371/journal.pone.0110268.
- Mai U, Mirarab S. 2018. TreeShrink: fast and accurate detection of outlier long branches in collections of phylogenetic trees. *BMC Genomics* 19(Suppl. 5):272 DOI 10.1186/s12864-018-4620-2.
- Mallet J, Besansky N, Hahn MW. 2015. How reticulated are species? *BioEssays* 38(2):140–149 DOI 10.1002/bies.201500149.
- Manceau M, Lambert A, Morlon H. 2015. Phylogenies support out-of-equilibrium models of biodiversity. *Ecology Letters* 18(4):347–356 DOI 10.1111/ele.12415.
- Manceau M, Lambert A, Morlon H. 2017. A unifying comparative phylogenetic framework including traits coevolving across interacting lineages. *Systematic Biology* 66:551–568.
- Mandáková T, Joly S, Krzywinski M, Mummenhoff K, Lysak MA. 2010. Fast diploidization in close mesopolyploid relatives of *Arabidopsis*. *Plant Cell* 22(7):2277–2290 DOI 10.1105/tpc.110.074526.
- Marcovitz A, Jia R, Bejerano G. 2016. “Reverse genomics” predicts function of human conserved noncoding elements. *Molecular Biology and Evolution* 33(5):1358–1369 DOI 10.1093/molbev/msw001.
- Marcussen T, Heier L, Brysting AK, Oxelman B, Jakobsen KS. 2015. From gene trees to a dated allopolyploid network: Insights from the angiosperm genus *Viola* (Violaceae). *Systematic Biology* 64(1):84–101 DOI 10.1093/sysbio/syu071.
- Marcussen T, Sandve SR, Heier L, Spannagl M, Pfeifer M, The International Wheat Genome Sequencing Consortium, Wulff BBH, Steuernagel B, Mayer KFX, Olsen O-A, Rogers J, Dole el J, Pozniak C, Eversole K, Feuillet C, Gill B, Friebe B, Lukaszewski AJ, Sourdille P, Endo TR, Kubalakova M, Ihalikova J, Dubska Z, Vrana J, Perkova R, Imkova H, Febrer M, Clissold L, McLay K, Singh K, Chhuneja P, Singh NK, Khurana J, Akhunov E, Choulet F, Alberti A, Barbe V, Wincker P, Kanamori H, Kobayashi F, Itoh T, Matsumoto T, Sakai H, Tanaka T, Wu J, Ogihara Y, Handa H, Maclachlan PR, Sharpe A, Klassen D, Edwards D, Batley J, Lien S, Caccamo M, Ayling S, Ramirez-Gonzalez RH, Clavijo BJ, Wright J, Martis MM, Mascher M, Chapman J, Poland JA, Scholz U, Barry K, Waugh R, Rokhsar DS, Muehlbauer GJ, Stein N, Gundlach H, Zytnicki M, Jamilloux V, Quesneville H, Wicker T, Faccioli P, Colaiacovo M, Stanca AM, Budak H, Cattivelli L, Glover N, Pingault L,

- Paux E, Sharma S, Appels R, Bellgard M, Chapman B, Nussbaumer T, Bader KC, Rimbart H, Wang S, Knox R, Kilian A, Alaux M, Alfama F, Couderc L, Guilhot N, Viseux C, Loaec M, Keller B, Praud S. 2014. Ancient hybridizations among the ancestral genomes of bread wheat. *Science* 345(6194):1250092 DOI 10.1126/science.1250092.
- Matzke A, Churakov G, Berkes P, Arms EM, Kelsey D, Brosius J, Kriegs JO, Schmitz J. 2012. Retroposon insertion patterns of neoavian birds: strong evidence for an extensive incomplete lineage sorting era. *Molecular Biology and Evolution* 29(6):1497–1501 DOI 10.1093/molbev/msr319.
- Matzke NJ, Wright A. 2016. Inferring node dates from tip dates in fossil Canidae: the importance of tree priors. *Biology Letters* 12(8):20160328 DOI 10.1098/rsbl.2016.0328.
- McCormack JE, Hird SM, Zellmer AJ, Carstens BC, Brumfield RT. 2013. Applications of next-generation sequencing to phylogeography and phylogenetics. *Molecular Phylogenetics and Evolution* 66(2):526–538 DOI 10.1016/j.ympev.2011.12.007.
- McCormack JE, Rodríguez-Gómez F, Tsai WLE, Faircloth BC, Webster MS. 2017. Transforming museum specimens into genomic resources. In: Webster MS, ed. *The Extended Specimen: Emerging Frontiers in Collections-Based Ornithological Research*. *Studies in Avian Biology* (no. 50). Boca Raton: CRC Press, 143–156.
- McCormack JE, Tsai WLE, Faircloth BC. 2016. Sequence capture of ultraconserved elements from bird museum specimens. *Molecular Ecology Resources* 16(5):1189–1203 DOI 10.1111/1755-0998.12466.
- McTavish EJ, Drew BT, Redelings B, Cranston KA. 2017. How and why to build a unified tree of life. *BioEssays* 39(11):1700114 DOI 10.1002/bies.201700114.
- McTavish EJ, Hinchliff CE, Allman JF, Brown J, Cranston KA, Holder MT, Rees JA, Smith SA. 2015. Phylsystem: a git-based data store for community-curated phylogenetic estimates: Fig. 1. *Bioinformatics* 31(17):2794–2800 DOI 10.1093/bioinformatics/btv276.
- Mendes FK, Fuentes-González JA, Schraiber JG, Hahn MW. 2018. A multispecies coalescent model for quantitative traits. *eLife* 7:e36482 DOI 10.7554/eLife.36482.
- Mendes FK, Hahn Y, Hahn MW. 2016. Gene tree discordance can generate patterns of diminishing convergence over time. *Molecular Biology and Evolution* 33(12):3299–3307 DOI 10.1093/molbev/msw197.
- Merker JD, Wenger AM, Sneddon T, Grove M, Zappala Z, Fresard L, Waggott D, Utiramerur S, Hou Y, Smith KS, Montgomery SB, Wheeler M, Buchan JG, Lambert CC, Eng KS, Hickey L, Korlach J, Ford J, Ashley EA. 2018. Long-read genome sequencing identifies causal structural variation in a Mendelian disease. *Genetics in Medicine* 20(1):159–163 DOI 10.1038/gim.2017.86.
- Meyer M, Kircher M, Gansauge M-T, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prufer K, de Filippo C, Sudmant PH, Alkan C, Fu Q, Do R, Rohland N, Tandon A, Siebauer M, Green RE, Bryc K, Briggs AW, Stenzel U, Dabney J, Shendure J, Kitzman J, Hammer MF, Shunkov MV, Derevianko AP, Patterson N, Andres AM, Eichler EE, Slatkin M, Reich D, Kelso J, Paabo S. 2012. A high-coverage genome sequence from an archaic Denisovan individual. *Science* 338(6104):222–226 DOI 10.1126/science.1224344.
- Mirarab S, Bayzid MS, Warnow T. 2016. Evaluating summary methods for multilocus species tree estimation in the presence of incomplete lineage sorting. *Systematic Biology* 65(3):366–380 DOI 10.1093/sysbio/syu063.
- Mirarab S, Warnow T. 2015. ASTRAL-II: coalescent-based species tree estimation with many hundreds of taxa and thousands of genes. *Bioinformatics* 31(12):i44–i52 DOI 10.1093/bioinformatics/btv234.

- Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, Frandsen PB, Ware J, Flouri T, Beutel RG, Niehuis O, Petersen M, Izquierdo-Carrasco F, Wappler T, Rust J, Aberer AJ, Aspöck U, Aspöck H, Bartel D, Blanke A, Berger S, Böhm A, Buckley TR, Calcott B, Chen J, Friedrich F, Fukui M, Fujita M, Greve C, Grobe P, Gu S, Huang Y, Jermiin LS, Kawahara AY, Krogmann L, Kubiak M, Lanfear R, Letsch H, Li Y, Li Z, Li J, Lu H, Machida R, Mashimo Y, Kapli P, McKenna DD, Meng G, Nakagaki Y, Navarrete-Heredia JL, Ott M, Ou Y, Pass G, Podsiadlowski L, Pohl H, von Reumont BM, Schütte K, Sekiya K, Shimizu S, Slipinski A, Stamatakis A, Song W, Su X, Szucsich NU, Tan M, Tan X, Tang M, Tang J, Timelthaler G, Tomizuka S. 2014. Phylogenomics resolves the timing and pattern of insect evolution. *Science* 346(6210):763–767 DOI 10.1126/science.1257570.
- Mitchell A, Mitter C, Regier JC. 2000. More taxa or more characters revisited: combining data from nuclear protein-encoding genes for phylogenetic analyses of Noctuoidea (Insecta: Lepidoptera). *Systematic Biology* 49(2):202–224 DOI 10.1093/sysbio/49.2.202.
- Montague MJ, Li G, Gandolfi B, Khan R, Aken BL, Searle SMJ, Minx P, Hillier LW, Koboldt DC, Davis BW, Driscoll CA, Barr CS, Blackstone K, Quilez J, Lorente-Galdos B, Marques-Bonet T, Alkan C, Thomas GWC, Hahn MW, Menotti-Raymond M, O'Brien SJ, Wilson RK, Lyons LA, Murphy WJ, Warren WC. 2014. Comparative analysis of the domestic cat genome reveals genetic signatures underlying feline biology and domestication. *Proceedings of the National Academy of Sciences of the United States of America* 111(48):17230–17235 DOI 10.1073/pnas.1410083111.
- Mora C, Tittensor DP, Adl S, Simpson AGB, Worm B. 2011. How many species are there on earth and in the ocean? *PLOS Biology* 9(8):e1001127 DOI 10.1371/journal.pbio.1001127.
- Morlon H. 2014. Phylogenetic approaches for studying diversification. *Ecology Letters* 17(4):508–525 DOI 10.1111/ele.12251.
- Morlon H, Parsons TL, Plotkin JB. 2011. Reconciling molecular phylogenies with the fossil record. *Proceedings of the National Academy of Sciences of the United States of America* 108(39):16327–16332 DOI 10.1073/pnas.1102543108.
- Murphy WJ, Larkin DM, Everts-van der Wind A, Bourque G, Tesler G, Auvil L, Beever JE, Chowdhary BP, Galibert F, Gatzke L, Hitte C, Meyers SN, Milan D, Ostrander EA, Pape G, Parker HG, Raudsepp T, Rogatcheva MB, Schook LB, Skow LC, Welge M, Womack JE, O'Brien SJ, Pevzner PA, Lewin HA. 2005. Dynamics of mammalian chromosome evolution inferred from multispecies comparative maps. *Science* 309(5734):613–617 DOI 10.1126/science.1111387.
- Murphy WJ, Pringle TH, Crider TA, Springer MS, Miller W. 2007. Using genomic data to unravel the root of the placental mammal phylogeny. *Genome Research* 17(4):413–421 DOI 10.1101/gr.5918807.
- Nabhan AR, Sarkar IN. 2012. The impact of taxon sampling on phylogenetic inference: a review of two decades of controversy. *Briefings in Bioinformatics* 13(1):122–134 DOI 10.1093/bib/bbr014.
- Nee S, Mooers AO, Harvey PH. 1992. Tempo and mode of evolution revealed from molecular phylogenies. *Proceedings of the National Academy of Sciences of the United States of America* 89(17):8322–8326 DOI 10.1073/pnas.89.17.8322.
- Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution* 32(1):268–274 DOI 10.1093/molbev/msu300.
- Ogilvie HA, Bouckaert RR, Drummond AJ. 2017. StarBEAST2 brings faster species tree inference and accurate estimates of substitution rates. *Molecular Biology and Evolution* 34(8):2101–2114 DOI 10.1093/molbev/msx126.

- Olave M, Sola E, Knowles LL. 2014. Upstream analyses create problems with DNA-based species delimitation. *Systematic Biology* 63(2):263–271 DOI 10.1093/sysbio/syt106.
- Oxelmann B, Yoshikawa N, McConaughy BL, Luo J, Denton AL, Hall BD. 2004. RPB2 gene phylogeny in flowering plants, with particular emphasis on asterids. *Molecular Phylogenetics and Evolution* 32(2):462–479 DOI 10.1016/j.ympev.2004.01.014.
- Pamilo P, Nei M. 1988. Relationships between gene trees and species trees. *Molecular Biology and Evolution* 5:568–583.
- Park SDE, Magee DA, McGettigan PA, Teasdale MD, Edwards CJ, Lohan AJ, Murphy A, Braud M, Donoghue MT, Liu Y, Chamberlain AT, Rue-Albrecht K, Schroeder S, Spillane C, Tai S, Bradley DG, Sonstegard TS, Loftus BJ, McHugh DE. 2015. Genome sequencing of the extinct Eurasian wild aurochs, *Bos primigenius*, illuminates the phylogeography and evolution of cattle. *Genome Biology* 16(1):234 DOI 10.1186/s13059-015-0790-2.
- Patel S, Kimball RT, Braun EL. 2013. Error in phylogenetic estimation for bushes in the tree of life. *Journal of Phylogenetics & Evolutionary Biology* 1(2):110 DOI 10.4172/2329-9002.1000110.
- Payseur BA, Rieseberg LH. 2016. A genomic perspective on hybridization and speciation. *Molecular Ecology* 25(11):2337–2360 DOI 10.1111/mec.13557.
- Pease JB, Haak DC, Hahn MW, Moyle LC. 2016. Phylogenomics reveals three sources of adaptive variation during a rapid radiation. *PLOS Biology* 14(2):e1002379 DOI 10.1371/journal.pbio.1002379.
- Pennell MW, FitzJohn RG, Cornwell WK. 2016. A simple approach for maximizing the overlap of phylogenetic and comparative data. *Methods in Ecology and Evolution* 7(6):751–758 DOI 10.1111/2041-210X.12517.
- Pennell MW, Harmon LJ. 2013. An integrative view of phylogenetic comparative methods: Connections to population genetics, community ecology, and paleobiology. *Annals of the New York Academy of Sciences* 1289:90–105.
- Peterson AT, Moyle RG, Nyári Á S, Robbins MB, Brumfield RT, Remsen JV Jr. 2007. The need for proper vouchering in phylogenetic studies of birds. *Molecular Phylogenetics and Evolution* 45(3):1042–1044 DOI 10.1016/j.ympev.2007.08.019.
- Pfeil BE, Brubaker CL, Craven LA, Crisp MD. 2004. Paralogy and orthology in the Malvaceae *rpb2* gene family: Investigation of gene duplication in *Hibiscus*. *Molecular Biology and Evolution* 21(7):1428–1437 DOI 10.1093/molbev/msh144.
- Philippe H, Brinkmann H, Lavrov DV, Littlewood DTJ, Manuel M, Wörheide G, Baurain D. 2011. Resolving difficult phylogenetic questions: why more sequences are not enough. *PLOS Biology* 9(3):e1000602 DOI 10.1371/journal.pbio.1000602.
- Philippe H, de Vienne DM, Ranwez V, Roure B, Baurain D, Delsuc F. 2017. Pitfalls in supermatrix phylogenomics. *European Journal of Taxonomy* 283:1–25.
- Piel WH, Chan L, Dominus MJ, Ruan J, Vos RA, Tannen V. 2009. Treebase v. 2: a database of phylogenetic knowledge. In: *e-BioSphere*. Available at <http://www.e-biosphere09.org>.
- Pleijel F, Jondelius U, Norlinder E, Nygren A, Oxelman B, Schander C, Sundberg P, Tholleson M. 2008. Phylogenies without roots? A plea for the use of vouchers in molecular phylogenetic studies. *Molecular Phylogenetics and Evolution* 48(1):369–371 DOI 10.1016/j.ympev.2008.03.024.
- Poe S. 1998. Sensitivity of phylogeny estimation to taxonomic sampling. *Systematic Biology* 47(1):18–31 DOI 10.1080/106351598261003.
- Potts AJ, Hedderson TA, Grimm GW. 2014. Constructing phylogenies in the presence of intra-individual site polymorphisms (2ISPs) with a focus on the nuclear ribosomal cistron. *Systematic Biology* 63(1):1–16 DOI 10.1093/sysbio/syt052.

- Prum RO, Berv JS, Dornburg A, Field DJ, Townsend JP. 2015. A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing. *Nature* 526(7574):569–573 DOI 10.1038/nature15697.
- Pyron RA. 2011. Divergence time estimation using fossils as terminal taxa and the origins of Lissamphibia. *Systematic Biology* 60(4):466–481 DOI 10.1093/sysbio/syr047.
- Pyron RA. 2015. Post-molecular systematics and the future of phylogenetics. *Trends in Ecology and Evolution* 30(7):384–389 DOI 10.1016/j.tree.2015.04.016.
- Pyron RA, Burbrink FT. 2012. Trait-dependent diversification and the impact of paleontological data on evolutionary hypothesis testing in New World ratsnakes (tribe Lampropeltini). *Journal of Evolutionary Biology* 25(3):497–508 DOI 10.1111/j.1420-9101.2011.02440.x.
- Quail MA, Smith M, Coupland P, Otto TD, Harris SR, Connor TR, Bertoni A, Swerdlow HP, Gu Y. 2012. A tale of three next generation sequencing platforms: comparison of Ion torrent, pacific biosciences and illumina MiSeq sequencers. *BMC Genomics* 13(1):341 DOI 10.1186/1471-2164-13-341.
- Quintero I, Keil P, Jetz W, Crawford FW. 2015. Historical biogeography using species geographical ranges. *Systematic Biology* 64(6):1059–1073 DOI 10.1093/sysbio/syv057.
- Rabosky DL. 2015. No substitute for real data: a cautionary note on the use of phylogenies from birth-death polytomy resolvers for downstream comparative analyses. *Evolution* 69(12):3207–3216 DOI 10.1111/evo.12817.
- Ramadugu C, Pfeil BE, Manjunath KL, Lee RF, Maureira-Butler IJ, Roose ML. 2013. A six nuclear gene phylogeny of citrus (Rutaceae) taking into account hybridization and lineage sorting. *PLOS ONE* 8(7):e68410 DOI 10.1371/journal.pone.0068410.
- Rannala B, Yang ZH. 2003. Bayes estimations of species divergence times and ancestral population sizes using DNA sequences from multiple loci. *Genetics* 164:1645–1656.
- Rannala B, Yang Z. 2008. Phylogenetic inference using whole genomes. *Annual Review of Genomics and Human Genetics* 9(1):217–231 DOI 10.1146/annurev.genom.9.081307.164407.
- Ranwez V, Delsuc F, Ranwez S, Belkhir K, Tilak M-K, Douzery EJP. 2007. OrthoMaM: a database of orthologous genomic markers for placental mammal phylogenetics. *BMC Evolutionary Biology* 7(1):241 DOI 10.1186/1471-2148-7-241.
- Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. 2014. Genome-wide inference of ancestral recombination graphs. *PLOS Genetics* 10(5):e1004342 DOI 10.1371/journal.pgen.1004342.
- Reddy S, Kimball RT, Pandey A, Hosner PA, Braun MJ, Hackett SJ, Han K-L, Harshman J, Huddleston CJ, Kingston S, Marks BD, Miglia KJ, Moore WS, Sheldon FH, Witt CC, Yuri T, Braun EL. 2017. Why do phylogenomic data sets yield conflicting trees? Data type influences the avian tree of life more than taxon sampling. *Systematic Biology* 66(5):857–879 DOI 10.1093/sysbio/syx041.
- Redon R, Ishikawa S, Fitch KR, Feuk L, Perry GH, Andrews TD, Fiegler H, Shapero MH, Carson AR, Chen W, Cho EK, Dallaire S, Freeman JL, González JR, Gratacòs MN, Huang J, Kalaitzopoulos D, Komura D, MacDonald JR, Marshall CR, Mei R, Montgomery L, Nishimura K, Okamura K, Shen F, Somerville MJ, Tchinda J, Valsesia A, Woodwark C, Yang F, Zhang J, Zerjal E. 2006. Global variation in copy number in the human genome. *Nature* 444(7118):444–454 DOI 10.1038/nature05329.
- Reid NM, Hird SM, Brown JM, Pelletier TA, McVay JD, Satler JD, Carstens BC. 2014. Poor fit to the multispecies coalescent is widely detectable in empirical data. *Systematic Biology* 63(3):322–333 DOI 10.1093/sysbio/syt057.
- Rhoads A, Au KF. 2015. PacBio sequencing and its applications. *Genomics, Proteomics and Bioinformatics* 13(5):278–289 DOI 10.1016/j.gpb.2015.08.002.

- Ricketts C, Popic V, Toosi H, Hajirasouliha I. 2018. Using LICHeE and BAMSE for reconstructing cancer phylogenetic trees. *Current Protocols in Bioinformatics* 62(1):e49 DOI 10.1002/cpbi.49.
- Rogozin IB, Thomson K, Csürös M, Carmel L, Koonin EV. 2008. Homoplasy in genome-wide analysis of rare amino acid replacements: the molecular-evolutionary basis for Vavilov's law of homologous series. *Biology Direct* 3(1):7 DOI 10.1186/1745-6150-3-7.
- Rogozin IB, Wolf YI, Babenko BN, Koonin EV. 2006. Dollo parsimony and the reconstruction of genome evolution. In: Albert VA, ed. *Parsimony, Phylogeny, and Genomics*. Oxford: Oxford University Press, 190–200.
- Rokas A, Carroll SB. 2005. More genes or more taxa? The relative contribution of gene number and taxon number to phylogenetic accuracy. *Molecular Biology and Evolution* 22(5):1337–1344 DOI 10.1093/molbev/msi121.
- Rokas A. 2011. Phylogenetic analysis of protein sequence data using the randomized accelerated maximum likelihood (RAXML) program. *Current Protocols in Molecular Biology* 96(1):19.11.1–19.11.14 DOI 10.1002/0471142727.mb1911s96.
- Rokas A, Holland PWH. 2000. Rare genomic changes as a tool for phylogenetics. *Trends in Ecology & Evolution* 15(11):454–459 DOI 10.1016/S0169-5347(00)01967-4.
- Rokas A, Williams BL, King N, Carroll SB. 2003. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature* 425(6960):798–804 DOI 10.1038/nature02053.
- Romiguier J, Cameron SA, Woodard SH, Fischman BJ, Keller L, Praz CJ. 2016. Phylogenomics controlling for base compositional bias reveals a single origin of eusociality in corbiculate bees. *Molecular Biology and Evolution* 33(3):670–678 DOI 10.1093/molbev/msv258.
- Romiguier J, Roux C. 2017. Analytical biases associated with GC-content in molecular evolution. *Frontiers in Genetics* 8(e72):16 DOI 10.3389/fgene.2017.00016.
- Roncal J, Guyot R, Hamon P, Crouzillat D, Rigoreau M, Konan ON, Rakotomalala JJ, Nowak MD, Davis AP, de Kochko A. 2016. Active transposable elements recover species boundaries and geographic structure in Madagascan coffee species. *Molecular Genetics and Genomics* 291(1):155–168 DOI 10.1007/s00438-015-1098-3.
- Ronquist F, Klopstein S, Vilhelmsen L, Schulmeister S, Murray DL, Rasnitsyn AP. 2012. A total-evidence approach to dating with fossils, applied to the early radiation of the Hymenoptera. *Systematic Biology* 61(6):973–999 DOI 10.1093/sysbio/sys058.
- Rosenberg NA, Nordborg M. 2002. Genealogical trees, coalescent theory and the analysis of genetic polymorphisms. *Nature Reviews Genetics* 3(5):380–390 DOI 10.1038/nrg795.
- Rosindell J, Cornell SJ, Hubbell SP, Etienne RS. 2010. Protracted speciation revitalizes the neutral theory of biodiversity. *Ecology Letters* 13(6):716–727 DOI 10.1111/j.1461-0248.2010.01463.x.
- Rosindell J, Harmon LJ. 2012. OneZoom: a fractal explorer for the tree of life. *PLOS Biology* 10(10):e1001406 DOI 10.1371/journal.pbio.1001406.
- Rosindell J, Harmon LJ, Etienne RS. 2015. Unifying ecology and macroevolution with individual-based theory. *Ecology Letters* 18:472–482 DOI 10.1111/ele.12430.
- Rota-Stabelli O, Campbell L, Brinkmann H, Edgecombe GD, Longhorn SJ, Peterson KJ, Pisani D, Philippe H, Telford MJ. 2011. A congruent solution to arthropod phylogeny: phylogenomics, microRNAs and morphology support monophyletic Mandibulata. *Proceedings of the Royal Society B: Biological Sciences* 278(1703):298–306 DOI 10.1098/rspb.2010.0590.
- Ruane S, Austin CC. 2017. Phylogenomics using formalin-fixed and 100+ year-old intractable natural history specimens. *Molecular Ecology Resources* 17(5):1003–1008 DOI 10.1111/1755-0998.12655.

- Ryan JF, Pang K, Schnitzler CE, Nguyen A-D, Moreland RT, Simmons DK, Koch BJ, Francis WR, Havlak P, Comparative Sequencing Program NISC, Smith SA, Putnam NH, Haddock SHD, Dunn CW, Wolfsberg TG, Mullikin JC, Martindale MQ, Baxevanis AD. 2013. The genome of the ctenophore *Mnemiopsis leidyi* and its implications for cell type evolution. *Science* 342(6164):1242592–1242592 DOI 10.1126/science.1242592.
- Sagitov S, Bartoszek K. 2012. Interspecies correlation for neutrally evolving traits. *Journal of Theoretical Biology* 309:11–19 DOI 10.1016/j.jtbi.2012.06.008.
- Sanderson MJ. 2002. Estimating absolute rates of molecular evolution and divergence times: a penalized likelihood approach. *Molecular Biology and Evolution* 19(1):101–109 DOI 10.1093/oxfordjournals.molbev.a003974.
- Sanderson MJ, Donoghue MJ, Piel WH, Eriksson T. 1994. TreeBASE: a prototype database of phylogenetic analyses and an interactive tool for browsing the phylogeny of life. *American Journal of Botany* 81(8):183.
- Sayyari E, Mirarab S. 2016. Fast coalescent-based computation of local branch support from quartet frequencies. *Molecular Biology and Evolution* 33(7):1654–1668 DOI 10.1093/molbev/msw079.
- Scally A, Dutheil JY, Hillier LDW, Jordan GE, Goodhead I, Herrero J, Hobolth A, Lappalainen T, Mailund T, Marques-Bonet T, McCarthy S, Montgomery SH, Schwalie PC, Tang YA, Ward MC, Xue Y, Yngvadottir B, Alkan C, Andersen LN, Ayub Q, Ball EV, Beal K, Bradley BJ, Chen Y, Clee CM, Fitzgerald S, Graves TA, Gu Y, Heath P, Heger A, Karakoc E, Kolb-Kokocinski A, Laird GK, Lunter G, Meader S, Mort M, Mullikin JC, Munch K, O'Connor TD, Phillips AD, Prado-Martinez J, Rogers AS, Sajjadian S, Schmidt D, Shaw K, Simpson JT, Stenson PD, Turner DJ, Vigilant L, Vilella AJ, Whitener W, Zhu B, Cooper DN, de Jong P, Dermitzakis ET, Eichler EE, Flicek P, Goldman N, Mundy NI, Ning Z, Odom DT, Ponting CP, Quail MA, Ryder OA, Searle SM, Warren WC, Wilson RK, Schierup MH, Rogers J, Tyler-Smith C, Durbin R. 2012. Insights into hominid evolution from the gorilla genome sequence. *Nature* 483(7388):169–175 DOI 10.1038/nature10842.
- Schrempf D, Minh BQ, De Maio N, von Haeseler A, Kosiol C. 2016. Reversible polymorphism-aware phylogenetic models and their application to tree inference. *Journal of Theoretical Biology* 407:362–370 DOI 10.1016/j.jtbi.2016.07.042.
- Schwartz R, Schäffer AA. 2017. The evolution of tumour phylogenetics: principles and practice. *Nature Reviews Genetics* 18(4):213–229 DOI 10.1038/nrg.2016.170.
- Scotch M, Sarkar IN, Mei C, Leaman R, Cheung K-H, Ortiz P, Singraur A, Gonzalez G. 2011. Enhancing phylogeography by improving geographical information from GenBank. *Journal of Biomedical Informatics* 44:S44–S47 DOI 10.1016/j.jbi.2011.06.005.
- Shen X-X, Hittinger CT, Rokas A. 2017. Contentious relationships in phylogenomic studies can be driven by a handful of genes. *Nature Ecology and Evolution* 1(5):0126 DOI 10.1038/s41559-017-0126.
- Shi C-M, Yang Z. 2018. Coalescent-based analyses of genomic sequence data provide a robust resolution of phylogenetic relationships among major groups of gibbons. *Molecular Biology and Evolution* 35(1):159–179 DOI 10.1093/molbev/msx277.
- Siepel A. 2009. Phylogenomics of primates and their ancestral populations. *Genome Research* 19(11):1929–1941 DOI 10.1101/gr.084228.108.
- Simion P, Philippe H, Baurain D, Jager M, Richter DJ, Di Franco A, Roure B, Satoh N, Quéinnec É, Ereskovsky A, Lapébie P, Corre E, Delsuc F, King N, Wörheide G, Manuel M. 2017. A large and consistent phylogenomic dataset supports sponges as the sister group to all other animals. *Current Biology* 27(7):958–967 DOI 10.1016/j.cub.2017.02.031.

- Siu-Ting K, Gower DJ, Pisani D, Kassahun R, Gebresenbet F, Menegon M, Mengistu AA, Saber SA, de Sá R, Wilkinson M, Loader SP. 2014. Evolutionary relationships of the critically endangered frog *Ericabatrachus baleensis* Largen, 1991 with notes on incorporating previously unsampled taxa into large-scale phylogenetic analyses. *BMC Evolutionary Biology* 14(1):44 DOI 10.1186/1471-2148-14-44.
- Sjödin P, Jakobsson M. 2012. Population genetic nature of copy number variation. Population Genetic Nature of Copy Number Variation. In: Feuk L, ed. *Genomic Structural Variants. Methods in Molecular Biology (Methods and Protocols)*. Vol. 838. New York: Springer.
- Slater GJ. 2015. Iterative adaptive radiations of fossil canids show no evidence for diversity-dependent trait evolution. *Proceedings of the National Academy of Sciences of the United States of America* 112(16):4897–4902 DOI 10.1073/pnas.1403666111.
- Smith SA, Brown JW. 2018. Constructing a broadly inclusive seed plant phylogeny. *American Journal of Botany* 105(3):302–314 DOI 10.1002/ajb2.1019.
- Smith SA, Moore MJ, Brown JW, Yang Y. 2015. Analysis of phylogenomic datasets reveals conflict, concordance, and gene duplications with examples from animals and plants. *BMC Evolutionary Biology* 15(1):150 DOI 10.1186/s12862-015-0423-0.
- Smith SA, O’Meara BC. 2012. treePL: divergence time estimation using penalized likelihood for large phylogenies. *Bioinformatics* 28(20):2689–2690 DOI 10.1093/bioinformatics/bts492.
- Solís-Lemus C, Ané C. 2016. Inferring phylogenetic networks with maximum pseudolikelihood under incomplete lineage sorting. *PLOS Genetics* 12(3):e1005896 DOI 10.1371/journal.pgen.1005896.
- Solís-Lemus C, Bastide P, Ané C. 2017. PhyloNetworks: a package for phylogenetic networks. *Molecular Biology and Evolution* 34(12):3292–3298 DOI 10.1093/molbev/msx235.
- Solís-Lemus C, Knowles LL, Ané C. 2015. Bayesian species delimitation combining multiple genes and traits in a unified framework. *Evolution* 69(2):492–507 DOI 10.1111/evo.12582.
- Song S, Liu L, Edwards SV, Wu S. 2012. Resolving conflict in eutherian mammal phylogeny using phylogenomics and the multispecies coalescent model. *Proceedings of the National Academy of Sciences of the United States of America* 112:E6079–E6079 DOI 10.1073/pnas.1518753112.
- Sousa F, Bertrand YJK, Doyle JJ, Oxelman B, Pfeil BE. 2017. Using genomic location and coalescent simulation to investigate gene tree discordance in *Medicago* L. *Systematic Biology* 66(6):934–949 DOI 10.1093/sysbio/syx035.
- Springer MS, Gatesy J. 2016. The gene tree delusion. *Molecular Phylogenetics and Evolution* 94:1–33 DOI 10.1016/j.ympev.2015.07.018.
- Staats M, Erkens RHJ, van de Vossenberg B, Wieringa JJ, Kraaijeveld K, Stielow B, Geml J, Richardson JE, Bakker FT. 2013. Genomic treasure troves: complete genome sequencing of herbarium and insect museum specimens. *PLOS ONE* 8(7):e69189 DOI 10.1371/journal.pone.0069189.
- Stadler T. 2009. On incomplete sampling under birth–death models and connections to the sampling-based coalescent. *Journal of Theoretical Biology* 261(1):58–66 DOI 10.1016/j.jtbi.2009.07.018.
- Stadler T. 2013. Recovering speciation and extinction dynamics based on phylogenies. *Journal of Evolutionary Biology* 26(6):1203–1219 DOI 10.1111/jeb.12139.
- Stadler T, Steel M. 2012. Distribution of branch lengths and phylogenetic diversity under homogeneous speciation models. *Journal of Theoretical Biology* 297:33–40 DOI 10.1016/j.jtbi.2011.11.019.
- Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30(9):1312–1313 DOI 10.1093/bioinformatics/btu033.

- Stephens M, Smith NJ, Donnelly P. 2001. A new statistical method for haplotype reconstruction from population data. *The American Journal of Human Genetics* 68(4):978–989 DOI 10.1086/319501.
- Stoltzfus A, O'Meara B, Whitacre J, Mounce R, Gillespie EL, Kumar S, Rosauer DF, Vos RA. 2012. Sharing and re-use of phylogenetic trees (and associated data) to facilitate synthesis. *BMC Research Notes* 5(1):574 DOI 10.1186/1756-0500-5-574.
- Struck TH. 2013. The impact of paralogy on phylogenomic studies—a case study on annelid relationships. *PLOS ONE* 8(5):e62892 DOI 10.1371/journal.pone.0062892.
- Suarez AV, Tutsui ND. 2004. The value of museum collections for research and society. *BioScience* 54(1):66–74.
- Suh A, Paus M, Kiefmann M, Churakov G, Franke FA, Brosius J, Kriegs JO, Schmitz J. 2011. Mesozoic retrotransposons reveal parrots as the closest living relatives of passerine birds. *Nature Communications* 2(1):443 DOI 10.1038/ncomms1448.
- Suh A, Smeds L, Ellegren H. 2015. The dynamics of incomplete lineage sorting across the ancient adaptive radiation of neoavian birds. *PLOS Biology* 13(8):e1002224 DOI 10.1371/journal.pbio.1002224.
- Sukumaran J, Knowles LL. 2017. Multispecies coalescent delimits structure, not species. *Proceedings of the National Academy of Sciences of the United States of America* 114(7):1607–1612 DOI 10.1073/pnas.1607921114.
- Talavera G, Castresana J. 2007. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Systematic Biology* 56(4):564–577 DOI 10.1080/10635150701472164.
- Tang J, Moret BME, Cui L, dePamphilis CW. 2004. Phylogenetic reconstruction from arbitrary gene-order data. In: *Proceedings. Fourth IEEE Symposium on Bioinformatics and Bioengineering*. Taiwan: IEEE, 592–599.
- 1000 Genomes Project ConsortiumGibbs RA, Boerwinkle E, Doddapaneni H, Han Y, Korchina V, Kovar C, Lee S, Muzny D, Reid JG, Zhu Y, Wang J, Chang Y, Feng Q, Fang X, Guo X, Jian M, Jiang H, Jin X, Lan T, Li G, Li J, Li Y, Liu S, Liu X, Lu Y, Ma X, Tang M, Wang B, Wang G, Wu H, Wu R, Xu X, Yin Y, Zhang D, Zhang W, Zhao J, Zhao M, Zheng X, Lander ES, Altshuler DM, Gabriel SB, Gupta N, Gharani N, Toji LH, Gerry NP, Resch AM, Flicek P, Barker J, Clarke L, Gil L, Hunt SE, Kelman G, Kulesha E, Leinonen R, McLaren WM, Radhakrishnan R, Roa A, Smirnov D, Smith RE, Streeter I, Thormann A, Toneva I, Vaughan B, Zheng-Bradley X, Bentley DR, Grocock R, Humphray S, James T, Kingsbury Z, Lehrach H, Sudbrak R, Albrecht MW, Amstislavskiy VS, Borodina TA, Lienhard M, Mertes F, Sultan M, Timmermann B, Yaspo M-L, Mardis ER, Wilson RK, Fulton L, Fulton R, Sherry ST, Ananiev V, Belaia Z, Beloslyudtsev D, Bouk N, Chen C, Church D, Cohen R, Cook C, Garner J, Hefferon T, Kimelman M, Liu C, Lopez J, Meric P, O'Sullivan C, Ostapchuk Y, Phan L, Ponomarov S, Schneider V, Shekhtman E, Sirotkin K, Slotta D, Zhang H, McVean GA, Durbin RM, Balasubramaniam S, Burton J, Danecsek P, Keane TM, Kolb-Kokocinski A, McCarthy S, Stalker J, Quail M, Schmidt JP, Davies CJ, Gollub J, Webster T, Wong B, Zhan Y, Auton A, Campbell CL, Kong Y, Marcketta A, Gibbs RA, Yu F, Antunes L, Bainbridge M, Muzny D, Sabo A, Huang Z, Wang J, Coin LJM, Fang L, Guo X, Jin X, Li G, Li Q, Li Y, Li Z, Lin H, Liu B, Luo R, Shao H, Xie Y, Ye C, Yu C, Zhang F, Zheng H, Zhu H, Alkan C, Dal E, Kahveci F, Marth GT, Garrison EP, Kural D, Lee W-P, Fung Leong W, Stromberg M, Ward AN, Wu J, Zhang M, Daly MJ, DePristo MA, Handsaker RE, Altshuler DM, Banks E, Bhatia G, del Angel G, Gabriel SB, Genovese G, Gupta N, Li H, Kashin S, Lander ES, McCarroll SA, Nemesh JC, Poplin RE, Yoon SC, Lihm J, Makarov V, Clark AG, Gottipati S, Keinan A, Rodriguez-Flores JL, Korbel JO, Rausch T, Fritz MH, Stütz AM, Flicek P, Beal K, Clarke L,

- Datta A, Herrero J, McLaren WM, Ritchie GRS, Smith RE, et al. 2015. A global reference for human genetic variation. *Nature* 526(7571):68–74 DOI 10.1038/nature15393.
- Todd EV, Black MA, Gemmell NJ. 2016. The power and promise of RNA-seq in ecology and evolution. *Molecular Ecology* 25(6):1224–1241 DOI 10.1111/mec.13526.
- Tonini JFR, Beard KH, Ferreira RB, Jetz W, Pyron RA. 2016. Fully-sampled phylogenies of squamates reveal evolutionary patterns in threat status. *Biological Conservation* 204:23–31 DOI 10.1016/j.biocon.2016.03.039.
- Toprak Z, Pfeil BE, Jones G, Marcussen T, Ertekin AS, Oxelman B. 2016. Species delimitation without prior knowledge: DISSECT reveals extensive cryptic speciation in the *Silene aegyptiaca* complex (Caryophyllaceae). *Molecular Phylogenetics and Evolution* 102:1–8 DOI 10.1016/j.ympev.2016.05.024.
- Troudet J, Vigness-Lebbe R, Grandcolas P, Legendre F. 2018. The increasing disconnection of primary biodiversity data from specimens: how does it happen and how to handle it? *Systematic Biology* 67(6):1110–1119 DOI 10.1093/sysbio/syy044.
- Turney S, Cameron ER, Cloutier CA, Buddle CM. 2015. Non-repeatable science: assessing the frequency of voucher specimen deposition reveals that most arthropod research cannot be verified. *PeerJ* 3(2):e1168 DOI 10.7717/peerj.1168.
- Ulitsky I. 2016. Evolution to the rescue: using comparative genomics to understand long non-coding RNAs. *Nature Reviews Genetics* 17(10):601–614 DOI 10.1038/nrg.2016.85.
- Urrutia E, Chen H, Zhou Z, Zhang NR, Jiang Y. 2018. Integrative pipeline for profiling DNA copy number and inferring tumor phylogeny. *Bioinformatics* 34(12):2126–2128 DOI 10.1093/bioinformatics/bty057.
- Villar D, Berthelot C, Alridge S, Rayner TF, Lukk M, Pignatelli M, Park TJ, Deaville R, Erichsen JT, Jasinska AJ, Turner JM, Bertelsen MF, Murchison EP, Flicek P, Odom DT. 2015. Enhancer evolution across 20 mammalian species. *Cell* 160(3):554–566 DOI 10.1016/j.cell.2015.01.006.
- Vision TJ. 2010. Open data and the social contract of scientific publishing. *BioScience* 60(5):330–331 DOI 10.1525/bio.2010.60.5.2.
- Warnow T. 2018. Supertree construction: opportunities and challenges [q-bio.PE]. Available at <http://arxiv.org/abs/1805.03530v1>.
- Wen D, Nakhleh L. 2018. Co-estimating reticulate phylogenies and gene trees from multi-locus sequence data. *Systematic Biology* 67(3):439–457 DOI 10.1093/sysbio/syx085.
- Wen D, Yu Y, Nakhleh L. 2016. Bayesian Inference of Reticulate Phylogenies under the Multispecies Network Coalescent. *PLOS Genetics* 12:e1006006 DOI 10.1371/journal.pgen.1006006.
- Wen D, Yu Y, Hahn MW, Nakhleh L. 2016. Reticulate evolutionary history and extensive introgression in mosquito species revealed by phylogenetic network analysis. *Molecular Ecology* 25(11):2361–2372 DOI 10.1111/mec.13544.
- Wesche PL, Gaffney DJ, Keightley PD. 2009. DNA sequence error rates in genbank records estimated using the mouse genome as a reference. *DNA Sequence* 15(5–6):362–364 DOI 10.1080/10425170400008972.
- Wiedenhoeft J, Brugel E, Schliep A. 2016. Fast bayesian inference of copy number variants using hidden Markov models with wavelet compression. *PLOS Computational Biology* 12(5):e1004871 DOI 10.1371/journal.pcbi.1004871.
- Will KP, Mishler BD, Wheeler QD. 2005. The perils of DNA Barcoding and the need for integrative taxonomy. *Systematic Biology* 54(5):844–851 DOI 10.1080/10635150500354878.

- Wilson G, Aruliah DA, Brown CT, Chue-Hong NP, Davis M, Guy RT, Haddock SHD, Huff KD, Mitchell IM, Plumbley MD, Waugh B, White EP, Wilson P. 2014. Best practices for scientific computing. *PLOS Biology* 12(1):e1001745 DOI 10.1371/journal.pbio.1001745.
- Xi Z, Liu L, Davis CC. 2015. Genes with minimal phylogenetic information are problematic for coalescent analyses when gene tree estimation is biased. *Molecular Phylogenetics and Evolution* 92:63–71 DOI 10.1016/j.ympev.2015.06.009.
- Xi Z, Ruhfel BR, Schaefer H, Amorim AM, Sugumaran M, Wurdack KJ, Endress PK, Matthews ML, Stevens PF, Mathews S, Davis CC. 2012. Phylogenomics and a posteriori data partitioning resolve the Cretaceous angiosperm radiation Malpighiales. *Proceedings of the National Academy of Sciences of the United States of America* 109(43):17519–17524 DOI 10.1073/pnas.1205818109.
- Xu B, Yang Z. 2016. Challenges in species tree estimation under the multispecies coalescent model. *Genetics* 204(4):1353–1368 DOI 10.1534/genetics.116.190173.
- Yang Z, Rannala B. 2014. Unguided species delimitation using DNA sequence data from multiple loci. *Molecular Biology and Evolution* 31(12):3125–3135 DOI 10.1093/molbev/msu279.
- Yu Y, Dong J, Liu KJ, Nakhleh L. 2014. Maximum likelihood inference of reticulate evolutionary histories. *Proceedings of the National Academy of Sciences of the United States of America* 111(46):16448–16453 DOI 10.1073/pnas.1407950111.
- Zanne AE, Tank DC, Cornwell WK, Eastman JM, Smith SA, FitzJohn RG, McGlinn DJ, O'Meara BC, Moles AT, Reich PB, Royer DL, Soltis DE, Stevens PF, Westoby M, Wright IJ, Aarssen L, Bertin RI, Calaminus A, Govaerts R, Hemmings F, Leishman MR, Oleksyn J, Soltis PS, Swenson NG, Warman L, Beaulieu JM. 2014. Three keys to the radiation of angiosperms into freezing environments. *Nature* 506(7486):89–92 DOI 10.1038/nature12872.
- Zhang G, Jarvis ED, Gilbert MTP. 2014. A flock of genomes. *Science* 346:1308–1309 DOI 10.1126/science.346.6215.1308.
- Zhang C, Ogilvie HA, Drummond AJ, Stadler T. 2018. Bayesian inference of species networks from multilocus sequence data. *Molecular Biology and Evolution* 35(2):504–517 DOI 10.1093/molbev/msx307.